

Reward Optimizing Recommendation with Deep Learning and Maximum Inner Product Search

Criteo AI Lab KDD Tutorial 2022

Imad Aouali, Amine Benhalloum, Martin Bompaire, Achraf Ait Sidi Hammou,
Sergey Ivanov, Benjamin Heymann, David Rohde, Otmane Sakhi, Flavian Vasile,
Maxime Vono

Introduction – The Practical Stuff

What are we trying to do ?

Recommend the **best items** to a **user**

- **Recommend:** Within the context of online advertising, low latency constraints, slates
- **Best:** Maximizing a given reward (clicks, sales, ...)
- **Items:** Tens of millions of products
- **User:** History of past user interactions, adapt to change in behaviour quickly, users enter and leave the system continually

Score millions of items very efficiently while adapting to changes in user behaviour

General framework

We are interested in decision rules of the form: $score(u, a) \sim u^T \beta_a$

u – "user embedding" of dimension D

β – The item embeddings of dimension $D \times P$

β_a – The embedding of item a

General framework

We are interested in decision rules of the form: $score(u, a) \sim \mathbf{u}^T \boldsymbol{\beta}_a$

Why ?

To find the best action we need:

$$a^* = \operatorname{argmax}_{a' \in \{1, \dots, P\}} score(u, a')$$

If P is large, then this is very slow for general $score$ functions.

Using Dot Product enables us to leverage (Approximate) Maximum Inner Product Search techniques.

General framework

We are interested in decision rules of the form: $score(u, a) \sim \mathbf{u}^T \boldsymbol{\beta}_a$

Why ?

What happens during training ? We can't evaluate each action, negative sampling techniques come in handy here.

General framework

We are interested in decision rules of the form: $score(x, a) \sim f_{\Xi}(x)^T \beta_a$

x – The “context” the information we use to personalize on.

f_{Ξ} - Mapping from “context” to embedding

General framework

We are interested in decision rules of the form: $score(x, a) \sim f_{\Xi}(x)^T \beta_a$

Why ?

- Users enter and leave the system constantly
- We want to adapt to change in behaviour
- Generalize across users

General framework

We are interested in decision rules of the form: $score(x, a) \sim f_{\Xi}(x)^T h_{\beta}(a)$

a – Set of initial pre-trained embeddings, product attributes and features ...

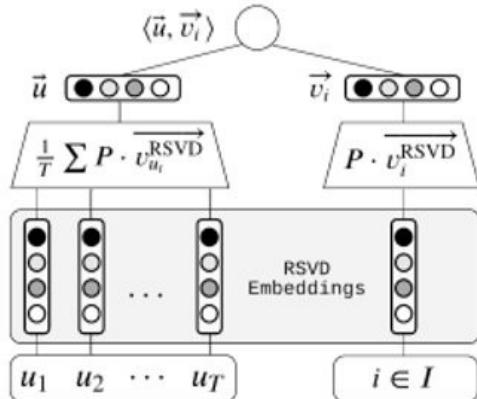
General framework

We are interested in decision rules of the form: $\textit{score}(\mathbf{x}, \mathbf{a}) \sim f_{\Xi}(\mathbf{x})^T h_{\beta}(\mathbf{a})$

Why ?

- Learn from vast amounts of heterogeneous sources of data, side information ...
- Enable transfer between tasks

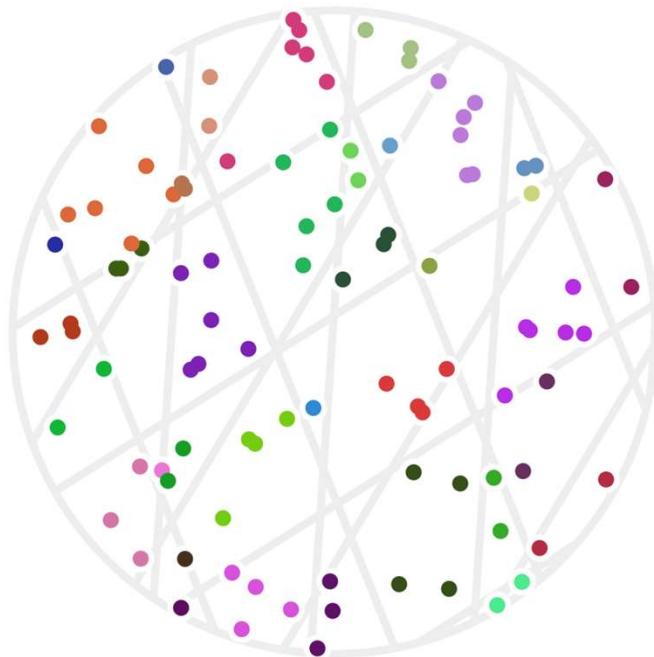
General framework



The now ubiquitous Two Tower style model.
([Criteo DeepKNN](#), [Facebook DLRM](#) ...)

You can make the towers as shallow or deep as you want/need #deeplearning

Maximum Inner Product Search



$\text{argsort}(f_{\Xi}(x)^T \beta)$

$$\begin{aligned} f_{\Xi}(x) &= D \\ \beta &= D \times P \end{aligned}$$

Sorting is $O(P \log P)$
Top-K is $O(P)$

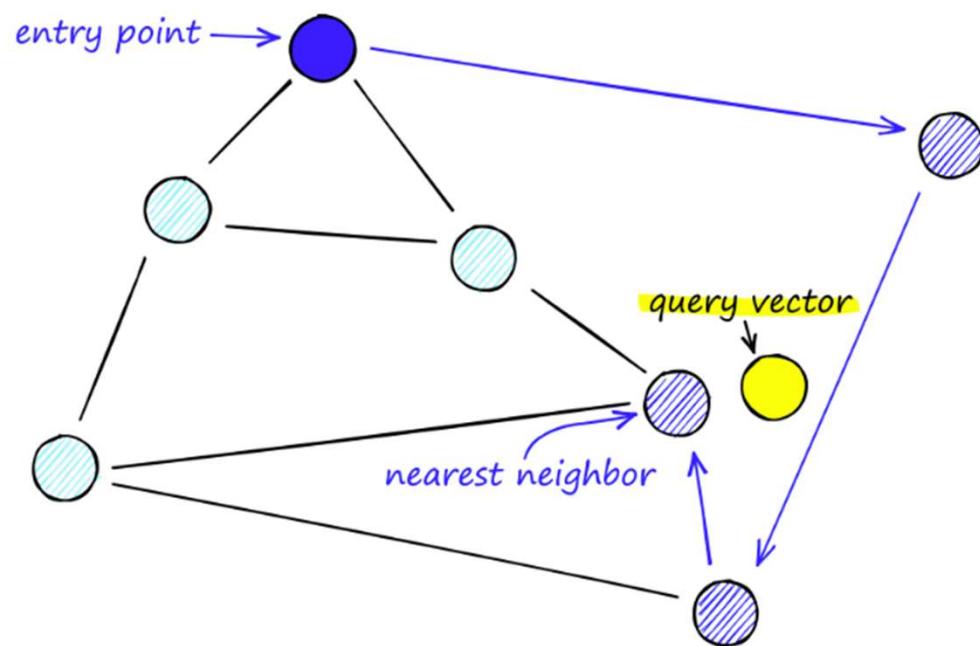
Approximate algorithms that “look like” $O(\log P)$

e.g. Navigable Small World

Side note: Hierarchical Navigable Small Worlds (HSNW)

General recipe for Graph Based NN methods

- Build a graph where close vectors are linked together
- At query time, route your query vector through the graph greedily, navigating to the closest neighbour at each time
- Stop when you hit a local minimum

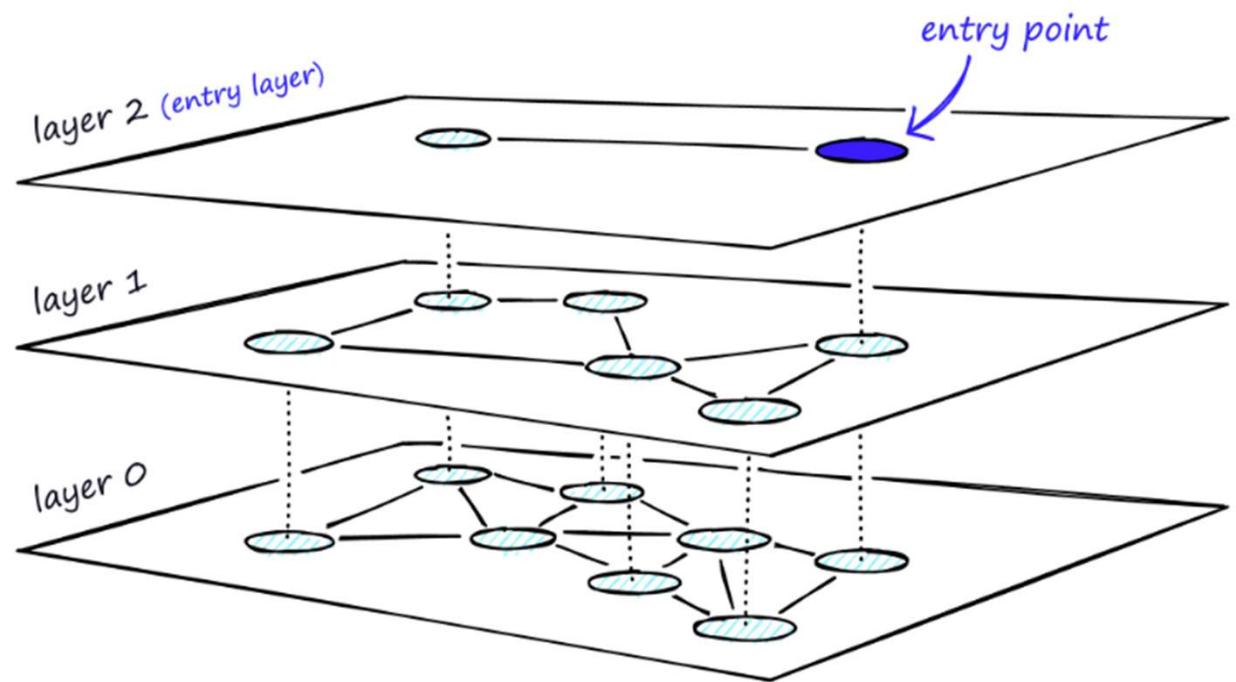


Credits: <https://www.pinecone.io/learn/hnsw/>

Side note: Hierarchical Navigable Small Worlds (HSNW)

The hierarchical part

- We start by a “zoomed out” view of the neighborhood graph
- And we zoom in after routing at each successive layer
- More fine grained search as we go deeper

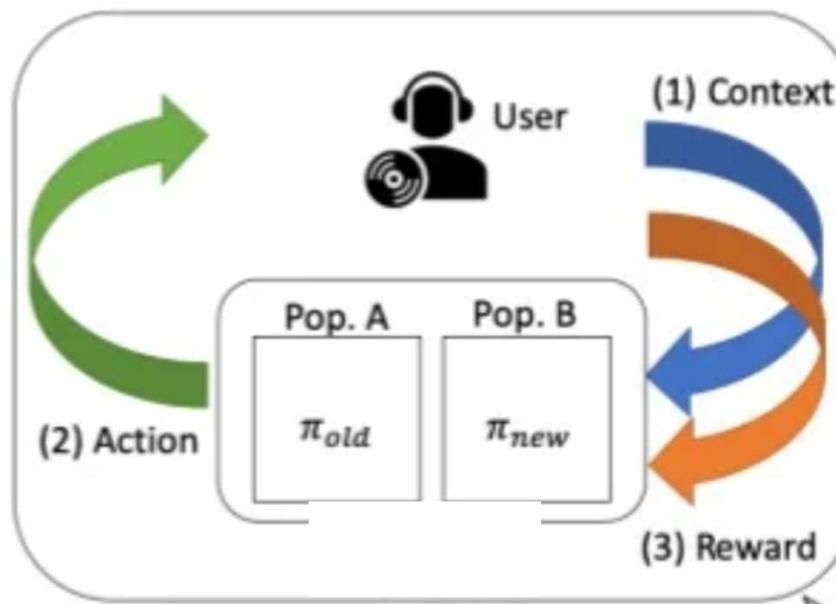


Credits: <https://www.pinecone.io/learn/hnsw/>

Introduction – The Science

Reward Optimizing Recommendation

Reward Optimizing Recommendation



Population A



Product view



Recommend slate



Product view



Recommend slate



No Sale



$$\text{slate} = A(\text{ } , \text{ })$$

$$\text{slate} = A(\text{ })$$

Population B



Product view



Product view



Recommend slate



Product view



Recommend slate



Sale



$$\text{slate} = B(\text{rice}, \text{rice}, \text{beer})$$

$$\text{slate} = B(\text{rice}, \text{rice})$$

Which is better A() or B()?

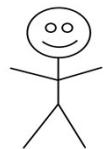
A() and B() are mappings from lists of items to lists of items e.g.

$$\begin{array}{c} \text{Beer} \\ \text{Phone} \end{array} = B(\text{Food 1}, \text{Food 2})$$

A() and B() may be constrained by engineering limitations

An A/B test can measure the timeline reward accurately (If we make the stable unit treatment value assumption).

How do we do this offline?



Product view



Product view



Recommend slate



Product view



Recommend slate



Sale



How do we do this offline?



Product view



Product view



Recommend slate



Product view



Recommend slate



Sale



This is really difficult to do!

Non-Reward Optimizing Recommendation with Pseudo Rewards

Our problem is that our recommender is a mapping:

$$\begin{array}{c} \text{Food} \\ \text{Beer} \\ \text{Smartphone} \end{array} = B(\text{Food}, \text{Beer})$$

And we're trying to answer the counterfactual question "What would happen if I had a different mapping ?" from data generated by a previous mapping, where rewards could be delayed, ...



Product
view





Product
view



Product
view





Product
view



Product
view



Product
view



Becomes a positive example





Product
view



Product
view



Product
view



$$L(\Xi, \boldsymbol{\beta}) = \log \sigma\{f_{\Xi}(v_1 = \text{red box}, v_2 = \text{brown box})^T \boldsymbol{\beta}_{a+} - f_{\Xi}(v_1 = \text{red box}, v_2 = \text{brown box})^T \boldsymbol{\beta}_{a-}\}$$

Becomes a positive example





Product view



Product view



Product view



$$L(\Xi, \beta) = \log \sigma\{f_{\Xi}(v_1 = \text{red box}, v_2 = \text{brown box})^T \beta_{a+} - f_{\Xi}(v_1 = \text{red box}, v_2 = \text{brown box})^T \beta_{a-}\}$$

Becomes a positive example

The negative might be
sampled uniformly or using
popularity

How do we use recommendation feedback ?

First answer, use a ranking loss function while training on historical recommendations and feedback.

e.g. Some pairwise ranking loss such as [BPR](#)



Product
view



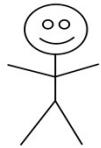


Product
view



Product
view





Product
view



Product
view



Recommend
Slate





Product view



Product view



Recommend Slate

No Click



Click



No Click



Product view



Product view



Recommend
Slate



$$L(\Xi, \beta) = \log \sigma\{f_{\Xi}(v_1 = \text{red box}, v_2 = \text{brown box})^T \beta_a^+ - f_{\Xi}(v_1 = \text{red box}, v_2 = \text{brown box})^T \beta_a^-\}$$

No Click

Click

No Click

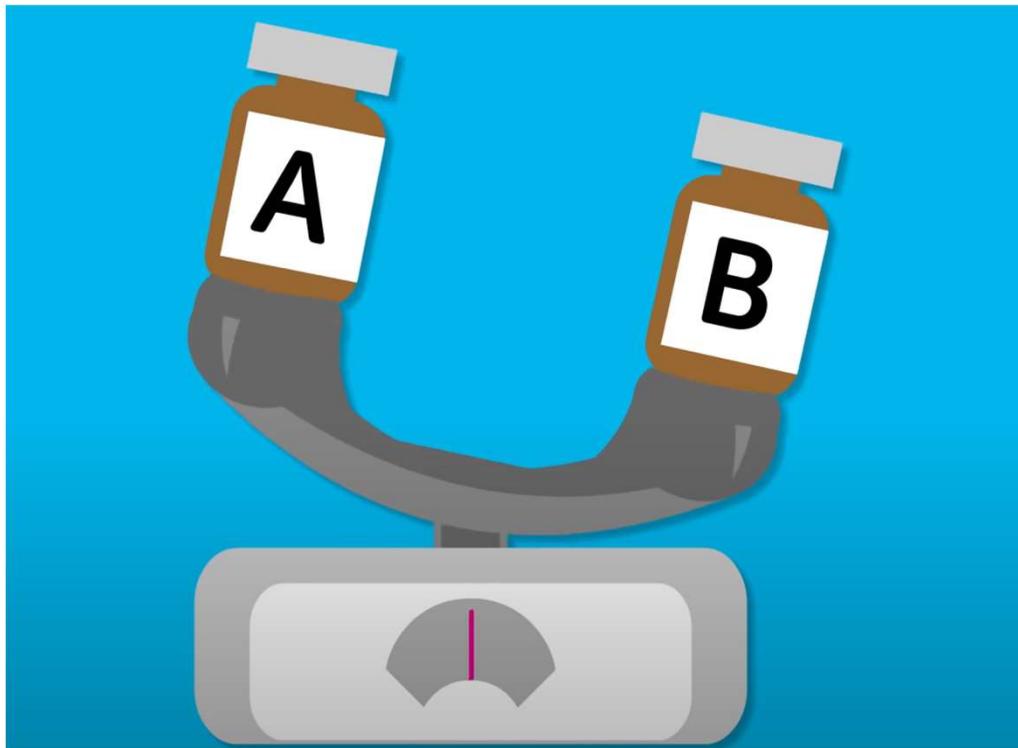
Notebook 1

[https://colab.research.google.com/github/otmhi/Reward-Optimizing-Reco/blob/main/Reward Optimizing Slate Recommendation with DL and MIPS Part 1.ipynb](https://colab.research.google.com/github/otmhi/Reward-Optimizing-Reco/blob/main/Reward%20Optimizing%20Slate%20Recommendation%20with%20DL%20and%20MIPS%20Part%201.ipynb)

Reward on the Timeline

Neither Evidence Based Medicine nor Reinforcement Learning offer easy off the shelf solutions

Randomized Control Trial in Medicine



- In medicine a randomized control trial, splits the population into two groups:
- We have:
 - A potentially strong intervention
 - Little personalization (clear membership)
 - Usually just two actions
 - Lots of negative results!



Product
View





Product
View



Product
View





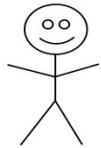
Product
View



Product
View



Recommend
Slate



Product
View



Product
View



Recommend
Slate

Slate =

$$\text{argsort}(f(v_1 = \text{[image of CousCous box]}, v_2 = \text{[image of Brown Basmati Rice bag]}))^T \boldsymbol{\beta})_{1:3}$$



Product
View



Product
View



Recommend
Slate



No click



Slate =

$$\text{argsort}(f(v_1 = \text{[image of CousCous box]}, v_2 = \text{[image of CousCous box]}))^T \boldsymbol{\beta})_{1:3}$$



Product
View



Product
View



Recommend
Slate



Product
View



Slate =

$$\text{argsort}(f(v_1 = \text{[image of CousCous box]}, v_2 = \text{[image of DAAWAT rice bag]}))^T \boldsymbol{\beta})_{1:3}$$



Product
View



Product
View



Recommend
Slate



Product
View



Recommend
Slate

$$\text{Slate} = \text{argsort}(f(v_1 = \text{CousCous}, v_2 = \text{DAAWAT Brown}, v_3 = \text{DAAWAT Almond}))^T \boldsymbol{\beta}_{1:2}$$

Slate =

$$\text{argsort}(f(v_1 = \text{CousCous}, v_2 = \text{DAAWAT Brown}, v_3 = \text{DAAWAT Almond}))^T \boldsymbol{\beta}_{1:3}$$



Product
View



Product
View



Recommend
Slate



Product
View



Recommend
Slate



Slate =

$$\text{argsort}(f(v_1 = \text{CousCous}, v_2 = \text{DAAWAT}, v_3 = \text{XXX}))^T \boldsymbol{\beta}_{1:2}$$

Slate =

$$\text{argsort}(f(v_1 = \text{CousCous}, v_2 = \text{DAAWAT}, v_3 = \text{XXX}))^T \boldsymbol{\beta}_{1:3}$$



Product
View



Product
View



Recommend
Slate



Product
View



Recommend
Slate



No click

Slate =

$$\text{argsort}(f(v_1 = \text{CousCous}, v_2 = \text{DAAWAT}, v_3 = \text{XXX}))^T \boldsymbol{\beta}_{1:2}$$

Slate =

$$\text{argsort}(f(v_1 = \text{CousCous}, v_2 = \text{DAAWAT}, v_3 = \text{XXX}))^T \boldsymbol{\beta}_{1:3}$$



Product
View



Product
View



Recommend
Slate



Product
View



Recommend
Slate



Sale



Slate =

$$\text{argsort}(f(v_1 = \text{CousCous}, v_2 = \text{DAAWAT}, v_3 = \text{XXX}))^T \boldsymbol{\beta}_{1:2}$$

Slate =

$$\text{argsort}(f(v_1 = \text{CousCous}, v_2 = \text{DAAWAT}, v_3 = \text{XXX}))^T \boldsymbol{\beta}_{1:3}$$

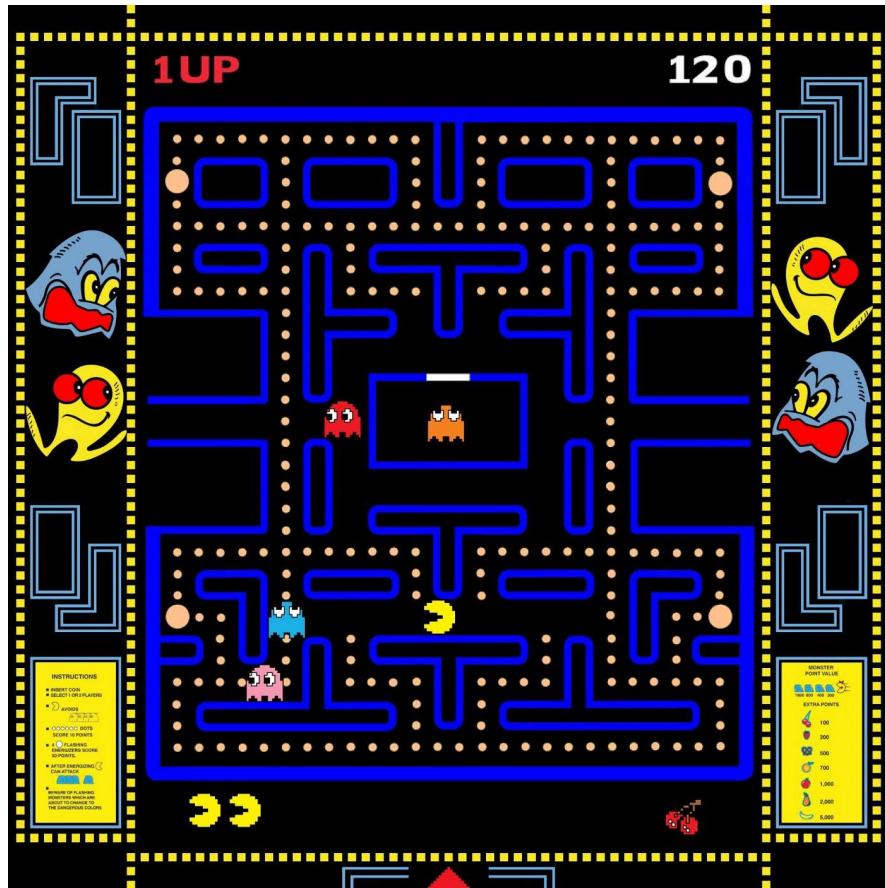
RCT in Medicine vs A/B Test for a Recommender System

	Medicine	RecSys
Personalization segments	<p>The user is who they are. Small number of categories e.g., age band or sex.</p>	<p>User is what they do. If they view an item – then that is who they are and what we want to personalize on.</p> <p>We learn about who the user is asynchronously.</p> <p>The interventions / recommendations might actually cause the events that we personalize on.</p>

RCT in Medicine vs A/B Test for a Recommender System

	Medicine	RecSys
Personalization segment size	Small: age bands, sex, simple medical history	Complex – everything we know.
Personalization action size	Often just intervention and placebo – (dose and dosing strategy)	With a single item slate – very big. With a large slate – the scale of the number of particles in the universe ¹
Effect size	Large, if it works.	Usually relatively small.

Reinforcement Learning for Game Playing



- In certain very successful RL problems the problem is both similar and different to recommendation
- We have:
 - No cost to losing many times before you win
 - Strong reward signals with little noise
 - Sequential steps of actions are critical to unlocking reward

Reinforcement Learning for RecSys



RecoGym

- RL for RecSys
- In principle it is the right formulation to move reward to the timeline
- We have:
 - High cost to performing poor actions in the wild
 - Lots of noise
 - Modest episodic features

Decision Theoretic Recommendation

Separate the Technology and the Science of Recommendation

Decision Theory Basics

What are the **states of nature** that might occur?

Rain / No Rain



Decision Theory Basics

What are the **actions** that we might do?

Carry an umbrella / Leave umbrella



Leonard “Jimmie” Savage in 1951

Act \ State	Rain	Shine
Carry	Inconvenience and wet feet	Inconvenience and slight embarrassment
Don't carry	Miserable drenching	Bliss unalloyed



Leonard “Jimmie” Savage in 1951

Act \ State	Rain	Shine
Carry	Inconvenience and wet feet	Inconvenience and slight embarrassment
Don't carry	Miserable drenching	Bliss unalloyed



Leonard “Jimmie” Savage in 1951

		Rain	Shine
Act	State		
Carry	Inconvenience and wet feet	Inconvenience and slight embarrassment	
Don't carry	Miserable drenching		Bliss unalloyed

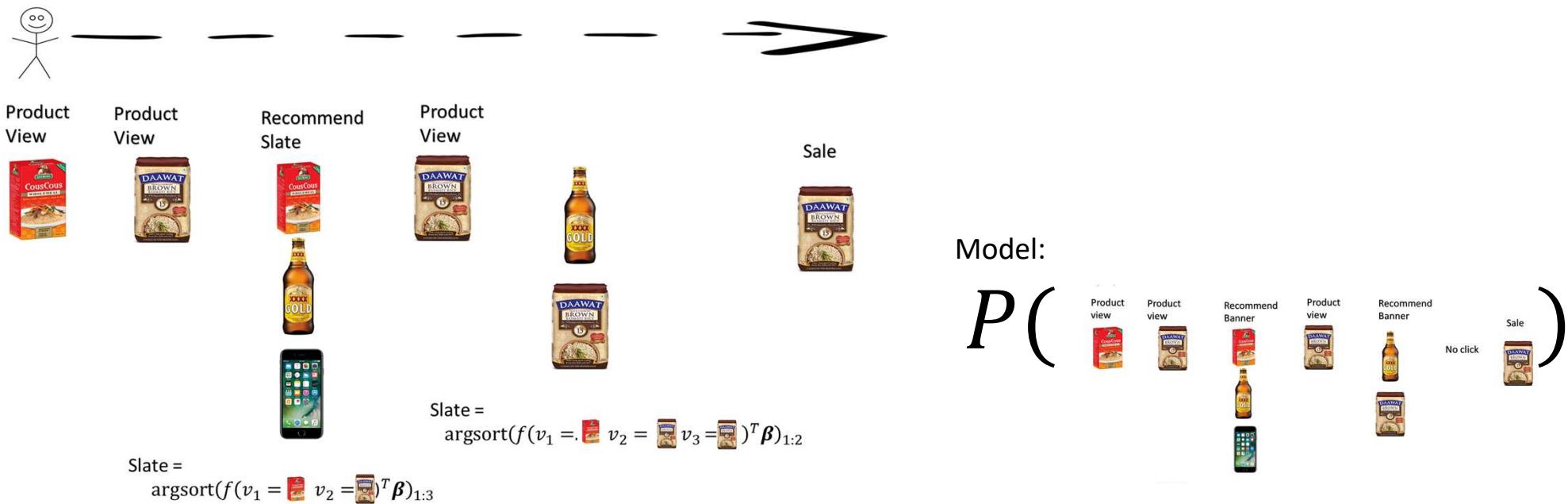
  

Leonard “Jimmie” Savage in 1951

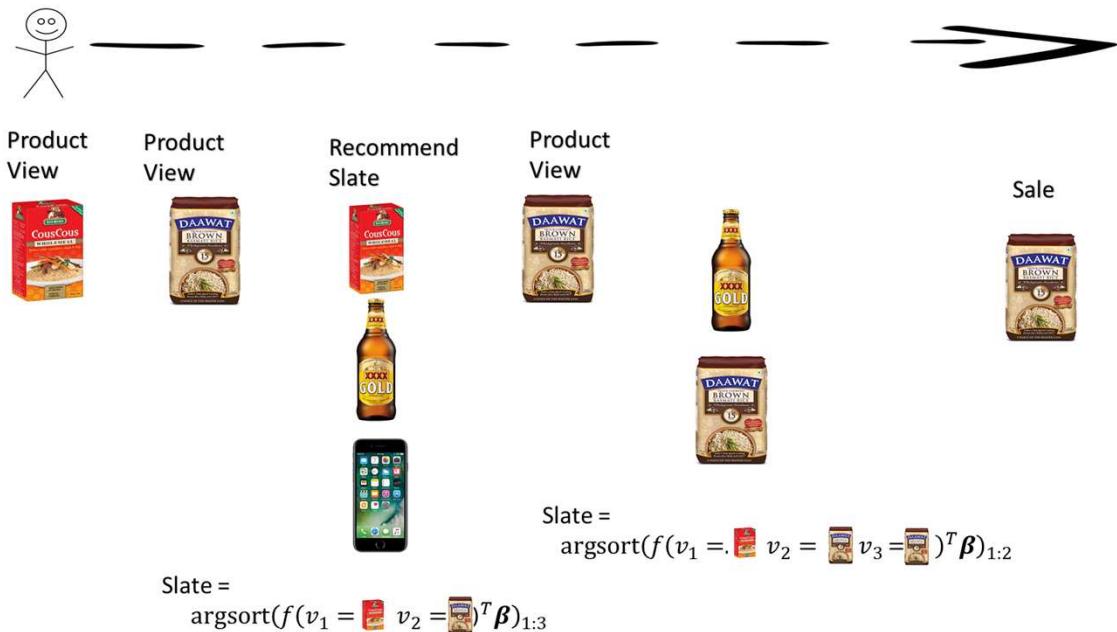
State		Rain	Shine
Act	Carry	Inconvenience and wet feet	Inconvenience and slight embarrassment
Don't carry	Miserable drenching		Bliss unalloyed

The model and the decision rule for RecSys



The model and the decision rule for RecSys



The model and the decision rule with the contextual bandit assumption:

Reward function (model):
 $R(a, x)$

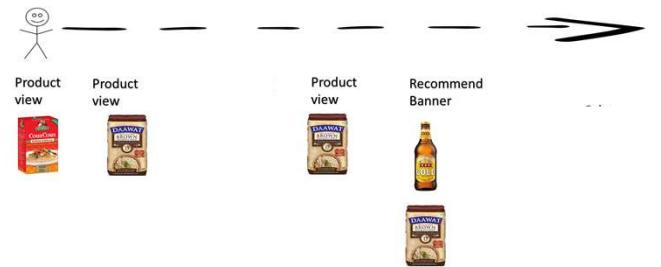
Marginal context (model)
 $P(x)$

Decision Rule
 $\text{argsort}(f(x)^T \boldsymbol{\beta})_{1:K}$

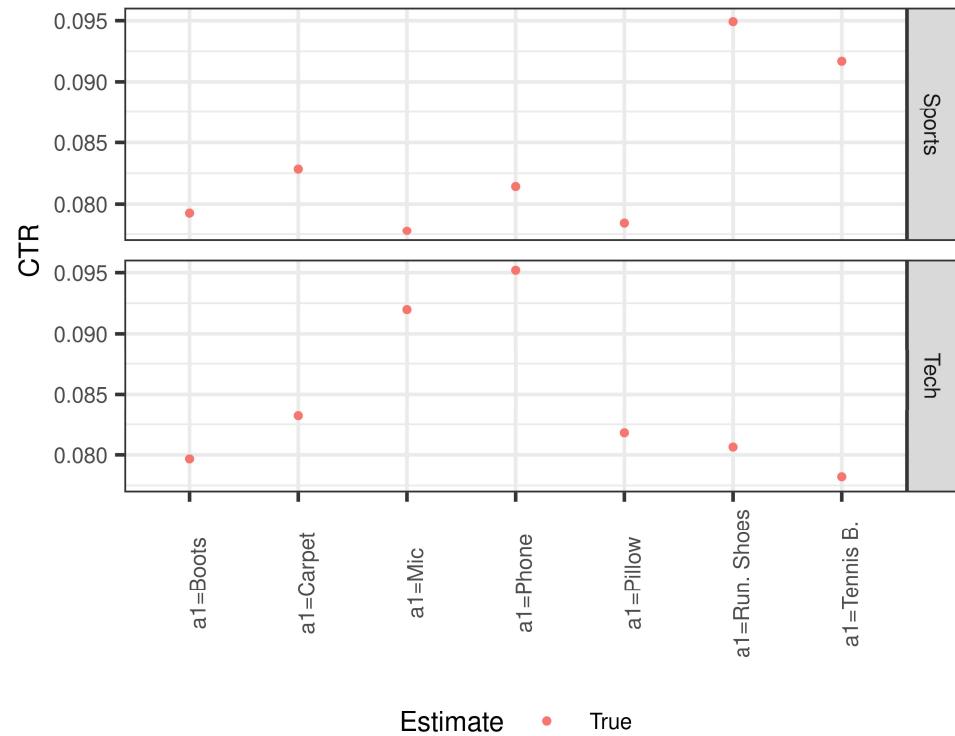
“Science”: Model based reward estimators

Estimate the reward of a recommendation using a model

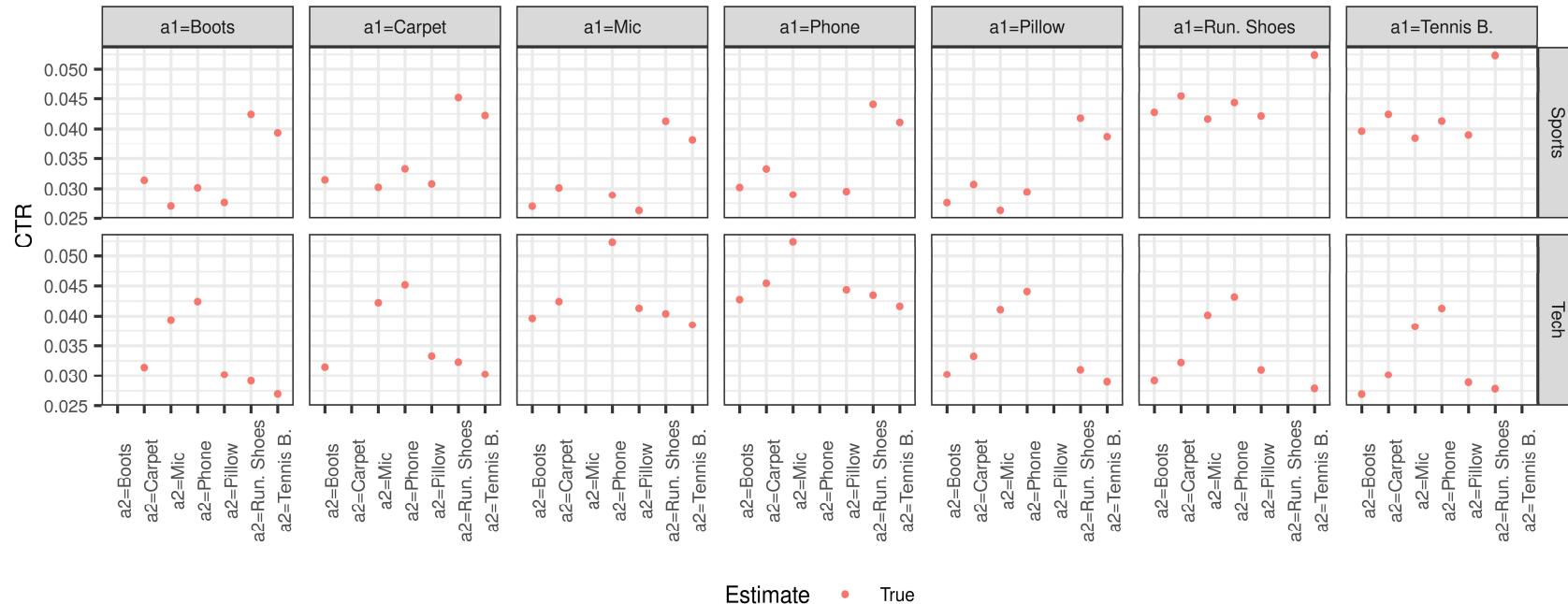
Contextual Bandit Assumption



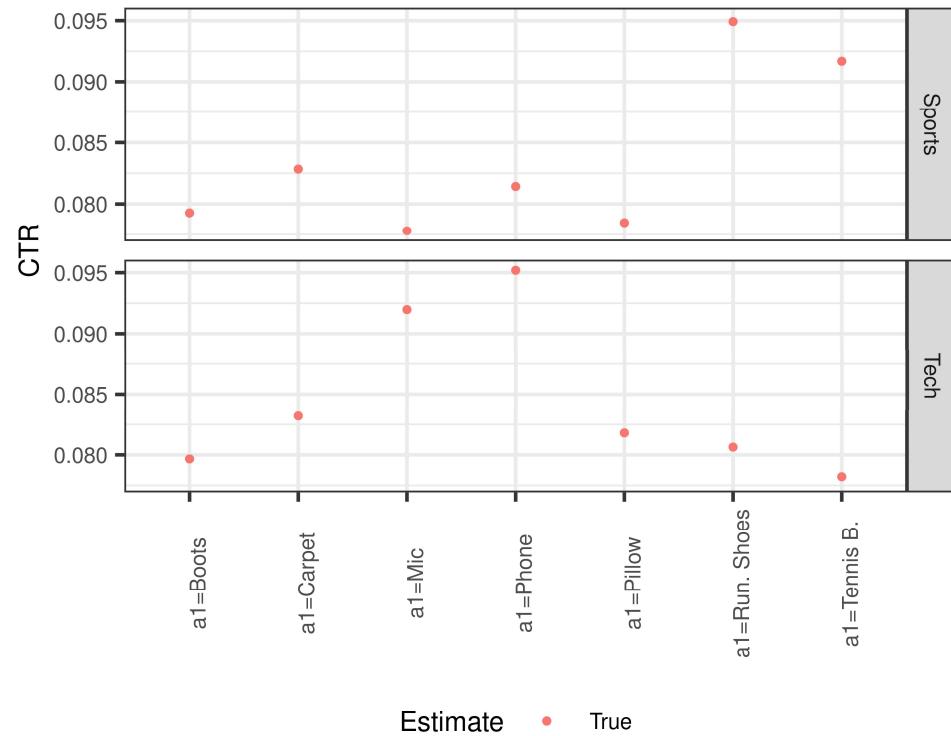
Oracle for a slate of size 1



Oracle for a slate of size 2

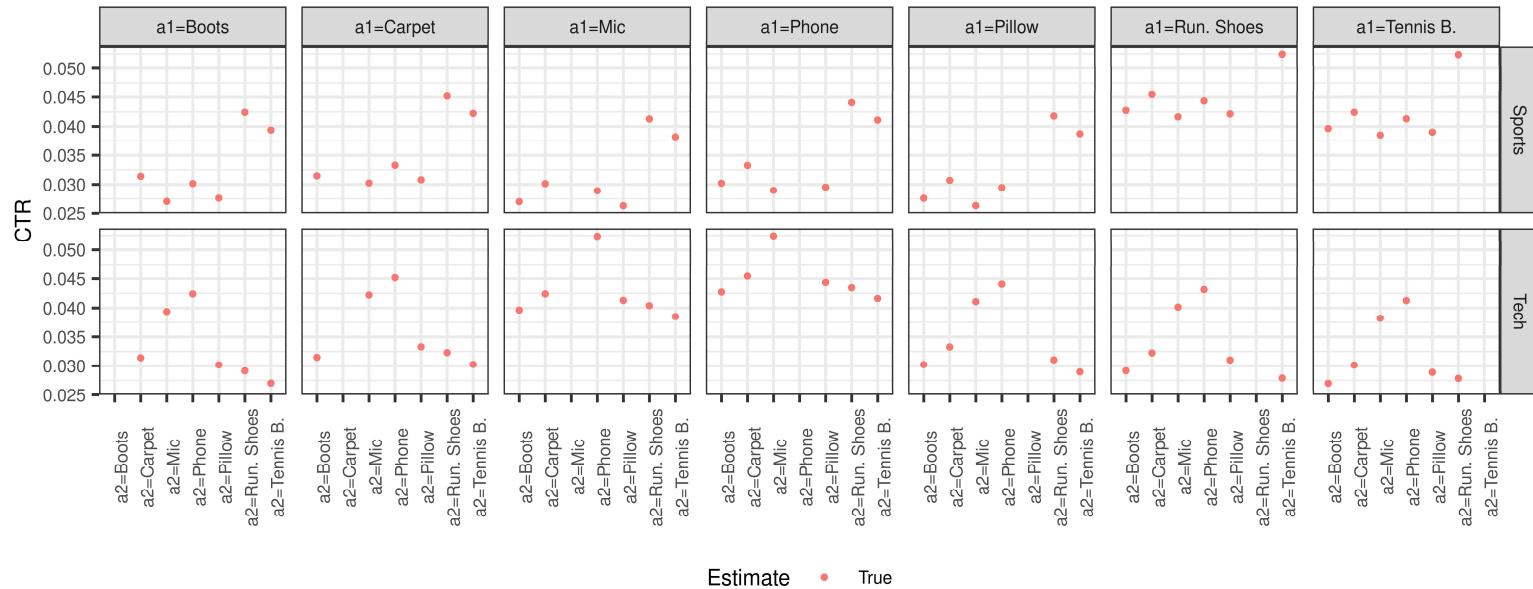


Oracle for a slate of size 1 – “Science”



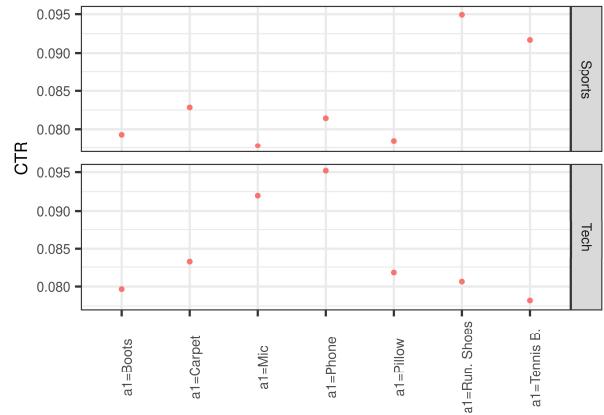
Estimate the CTRs
everywhere with a slate of
size 1

Oracle for a slate of size 2 – “Science”



Estimate the CTRs
everywhere with a slate of
size 2

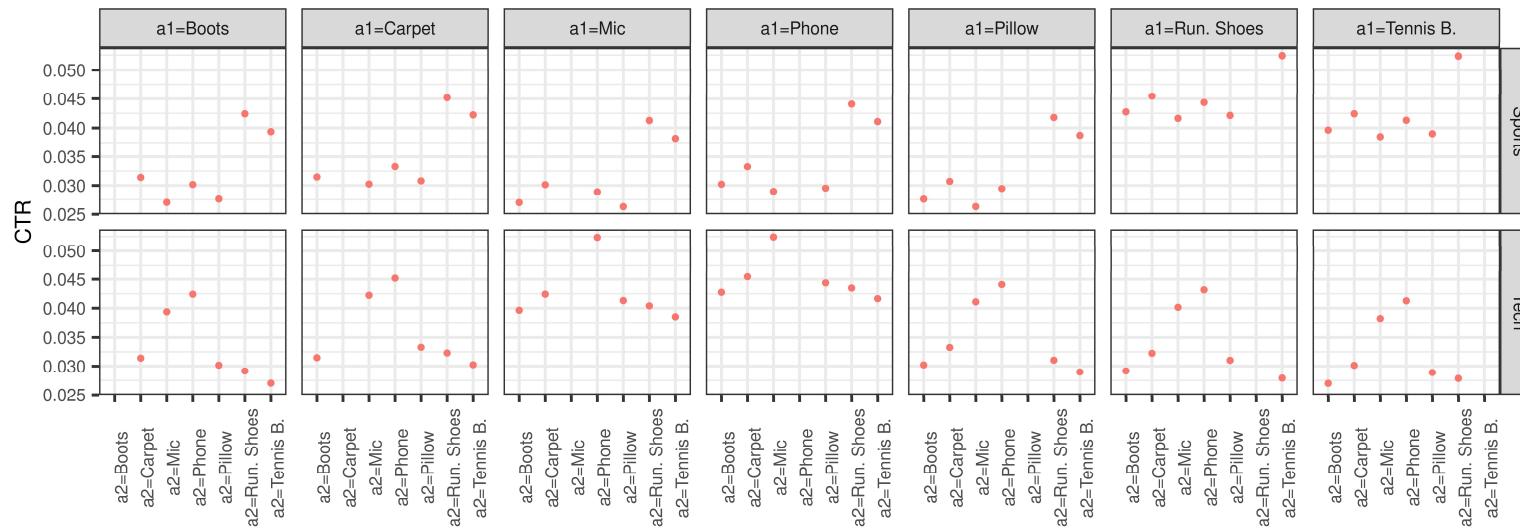
Oracle – “Science”



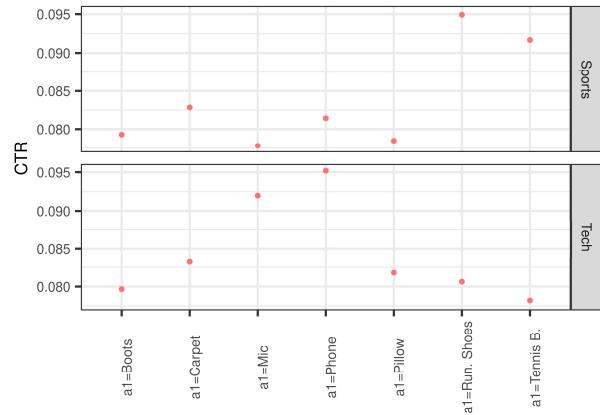
This is the reward function $R(a_1, x)$ and $R(a_1, a_2, x)$

... but we need to estimate $R(\cdot)$ from the recommender system log.

If we have logs of having tried all the recommendations with all the contexts – we can produce an estimate.

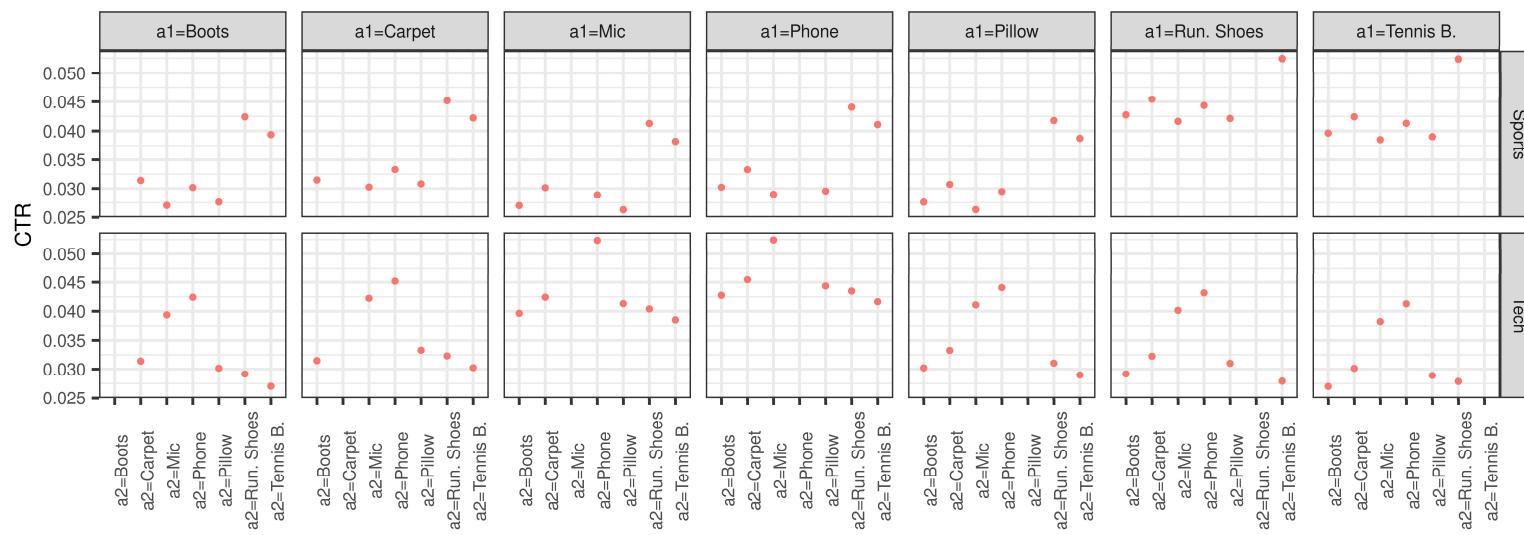


Oracle – “Science”

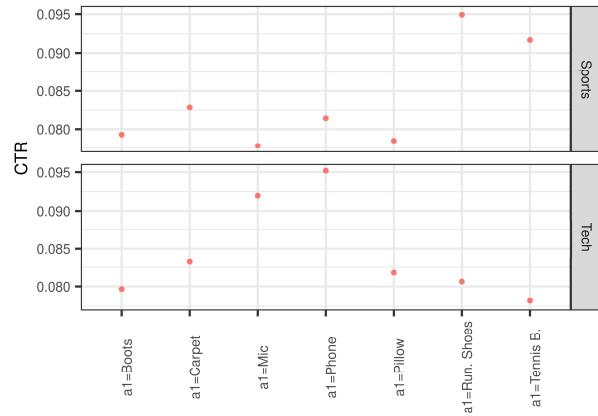


Three ways to estimate $R(\cdot)$:

- Apply directly to the real world
- Reward model
- Horvitz-Thompson Estimator

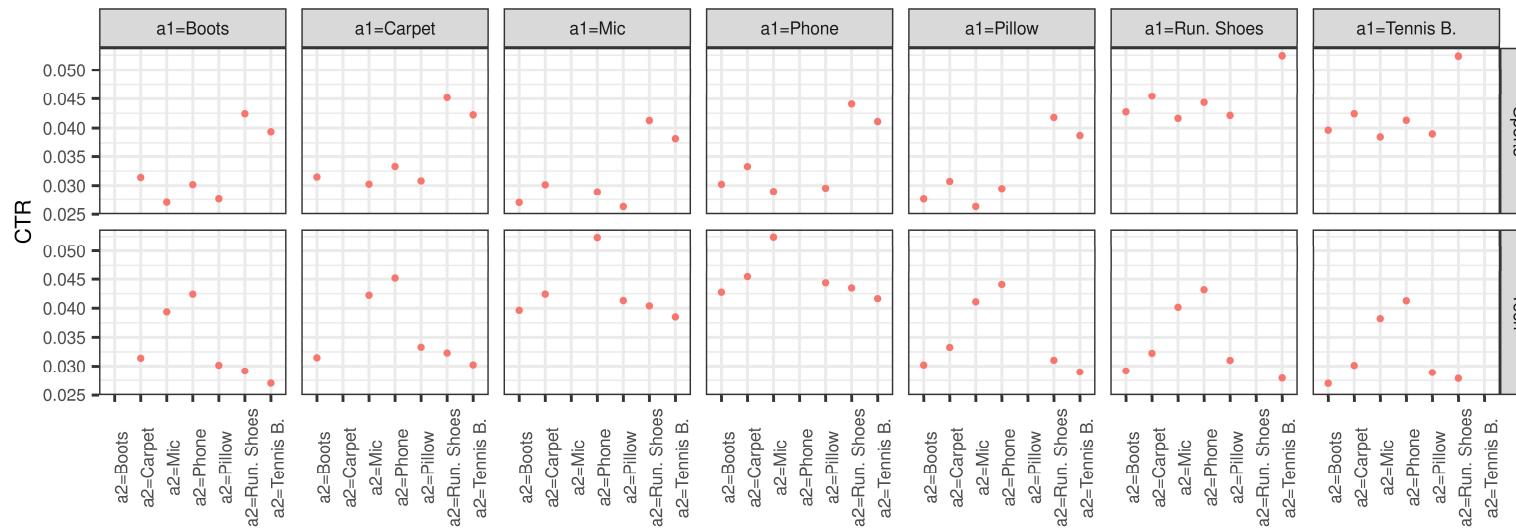


Oracle – “Technology”



$$\text{Run shoes} = \operatorname{argmax}(f_{\Xi}(\text{Sports})^T \boldsymbol{\beta})$$

$$\text{Mic} = \operatorname{argmax}(f_{\Xi}(\text{Tech})^T \boldsymbol{\beta})$$



$$\text{Run shoes, Ten Balls} = \operatorname{argsort}(f_{\Xi}(\text{Sports})^T \boldsymbol{\beta})_{1:2}$$

$$\text{Mic, Phone} = \operatorname{argsort}(f_{\Xi}(\text{Tech})^T \boldsymbol{\beta})_{1:2}$$

Single Item Recommendations

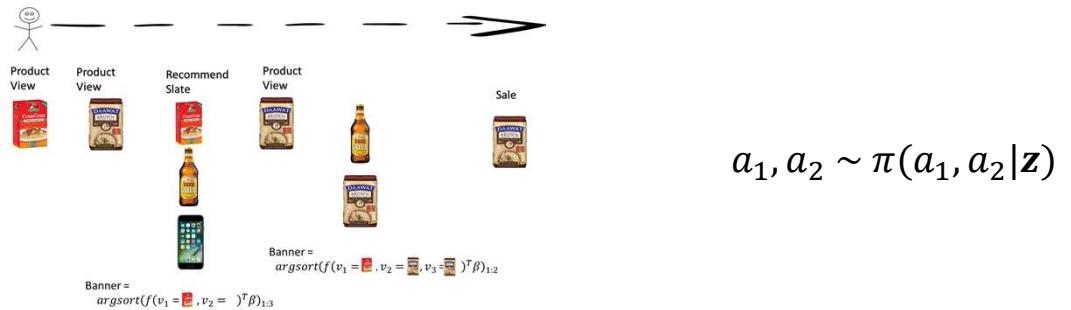
Just do it!

Just do it!

One way to get an unbiased estimate of $R(a, \mathbf{x})$ is just to do it!



- Unbiased
- You do not need to make the contextual bandit assumption



- You might go broke before you converge
- Noisy

Single Item Recommendations

Model Based Reward Estimation

Logistic Regression Reward Model (e.g. BLOB)

$$R(a, \mathbf{x}) = P(c = 1 | a, \mathbf{x}) = \sigma(g_\Gamma(\mathbf{x})^T \boldsymbol{\Psi}_a)$$

BLOB – used priors to accelerate learning the reward signal

Because of the parameterization of the reward model, we automatically get a MIPS:

$$a^* = \operatorname{argmax}_a (f_\Xi(\mathbf{z})^T \boldsymbol{\beta})$$

with $f_\Xi(\mathbf{x}) = g_\Gamma(\mathbf{x})$ and $\boldsymbol{\beta} = \boldsymbol{\Psi}$

Logistic Regression Reward Model (e.g. BLOB)

$$R(a, \mathbf{x}) = P(c = 1 | a, \mathbf{x}) = \sigma(g_{\Gamma}(\mathbf{x})^T \boldsymbol{\Psi})$$

BLOB – used priors to accelerate learning the reward signal

Because of the parameterization of the reward model, we automatically get a MIPS:

$$a^* = \operatorname{argmax}_a (f_{\Xi}(\mathbf{z})^T \boldsymbol{\beta})$$

with $f_{\Xi}(\mathbf{x}) = g_{\Gamma}(\mathbf{x})$ and $\boldsymbol{\beta} = \boldsymbol{\Psi}$



This is sufficient, but not necessary to build the MIPS, and it can result in larger than necessary embeddings

Logistic Regression Reward Model (e.g. BLOB)

$$R(a, \mathbf{x}) = P(c = 1 | a, \mathbf{x}) = \sigma(g_{\Gamma}(\mathbf{x})^T \boldsymbol{\Psi})$$

BLOB – used priors to accelerate learning the reward signal

Because of the parameterization of the reward model, we automatically get a MIPS:

$$a^* = \operatorname{argmax}_a (f_{\Xi}(\mathbf{z})^T \boldsymbol{\beta})$$

with $f_{\Xi}(\mathbf{x}) = g_{\Gamma}(\mathbf{x})$ and $\boldsymbol{\beta} = \boldsymbol{\Psi}$



We will discuss alternative ways to get $f_{\Xi}(\cdot), \boldsymbol{\beta}$

Logistic Regression Reward Model (e.g. BLOB)

We can split $x = [y, z]$

Where y - User Engagement, and z - User interest

$$R(a, x) = P(c = 1 | a, x) = \sigma(y^T \phi + g_\Gamma(z)^T \Psi)$$

$$a^* = \operatorname{argmax}_a (f_\Xi(z)^T \beta)$$

with $f_\Xi(z) = g_\Gamma(z)$ and $\beta = \Psi$

Logistic Regression Reward Model (e.g. BLOB)

We can split $x = [y, z]$

Where y - User Engagement, and z - User interest

$$R(a, x) = P(c = 1|a, x) = \sigma(y^T \phi + g_\Gamma(z)^T \Psi)$$

$$a^* = \operatorname{argmax}_a (f_\Xi(z)^T \beta)$$

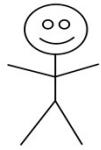
Nuisance parameters

with $f_\Xi(z) = g_\Gamma(z)$ and $\beta = \Psi$



Product
View





Product
View



Product
View





Product
View



Product
View



Recommend



Product
View



Product
View



Recommend



Click/No click

Recommend =
 $\text{argmax}(f(v_1 = \text{[cereal icon]}, v_2 = \text{[rice icon]}))^T \boldsymbol{\beta})$



Product
View



Product
View



Recommend



Click/No click

When we build a model to predict the reward, we can try to predict the reward (click).

Build a model based on:



We can also use features like:

- How active is the user?
- How big/attractive is the recommendation?

$$\text{Recommend} = \text{argmax}(f(v_1 = \text{CousCous}, v_2 = \text{Rice})^T \boldsymbol{\beta})$$



Product
View



Product
View



Recommend



Click/No click

This (often)
dominates in practice



Recommend =
 $\text{argmax}(f(v_1 = \text{[CousCous Box]}, v_2 = \text{[Rice Bag]})^T \boldsymbol{\beta})$

When we build a model to predict the reward, we can try to predict the reward (click).

Build a model based on:



We can also use features like:

- How active is the user?
- How big/attractive is the recommendation?



Product
View



Product
View



Recommend



Click/No click

We can predict
well with *only*
these features



Recommend =
 $\text{argmax}(f(v_1 = \text{[product image]}, v_2 = \text{[product image]}))^T \boldsymbol{\beta}$

When we build a model to predict the reward, we can try to predict the reward (click).

We can also use features like:

- How active is the user?
- How big/attractive is the recommendation?

Single Item Recommendations

Horvitz-Thompson Reward Estimation

Conditional Horvitz-Thompson Reward Estimation

$$R(a^*, \mathbf{x}) = \int c P(c|a', \mathbf{x}) \pi^*(a'|\mathbf{x}) da' = \int c \pi_0(a'|\mathbf{x}) \frac{P(c|a', \mathbf{x}) \pi^*(a'|\mathbf{x})}{\pi_0(a'|\mathbf{x})} da'$$

$\pi^*(a|\mathbf{x}) = 1$ if $a = a^*$ and $\pi^*(a|\mathbf{x}) = 0$, otherwise.

Horvitz-Thompson Reward Estimator

$$R(a, \mathbf{x}) = \frac{1}{N_{a,x}} \sum_{i \in (a_i=a, x_i=x)} \frac{c_i}{\pi_0(a_i | \mathbf{x}_i)}$$



We compute a sum over
the relevant history

Horvitz-Thompson Reward Estimator

Reward



$$R(a, \mathbf{x}) = \frac{1}{N_{a,x}} \sum_{i \in (a_i=a, x_i=x)} \frac{c_i}{\pi_0(a_i | \mathbf{x}_i)}$$

Horvitz-Thompson Reward Estimator

$$R(a, \mathbf{x}) = \frac{1}{N_{a,x}} \sum_{i \in (a_i=a, x_i=x)} \frac{c_i}{\pi_0(a_i | \mathbf{x}_i)}$$


Propensity score

Horvitz-Thompson Reward Estimator

$$R(a, \mathbf{x}) = \frac{1}{N_{a,x}} \sum_{i \in (a_i=a, x_i=x)} \frac{c_i}{\pi_0(a_i | \mathbf{x}_i)}$$



Number of records

Horvitz-Thompson Reward Estimator of a Policy

Previously we had a mapping $A(\mathbf{x})$ which returned a recommendation for user \mathbf{x} . This mapping is non-differentiable, so we introduce a “policy” $\pi(a|\mathbf{x})$. While $\pi(a|\mathbf{x})$ is mathematically a probability, but it does not represent uncertainty, it simply makes the (in this case argmax) differentiable. We then get:

$$E[U|\pi(\cdot)] = \frac{1}{N} \sum_i \frac{c_i \pi(a_i|\mathbf{x}_i)}{\pi_0(a_i|\mathbf{x}_i)}$$

Horvitz-Thompson Reward Estimator of a Policy

Previously we had a mapping $A(\mathbf{x})$ which returned a recommendation for user \mathbf{x} . This mapping is non-differentiable, so we introduce a “policy” $\pi(a|\mathbf{x})$. While $\pi(a|\mathbf{x})$ is mathematically a probability, but it does not represent uncertainty, it simply makes the (in this case argmax) differentiable. We then get:

$$E[U|\pi(\cdot)] = \frac{1}{N} \sum_i \frac{c_i \pi(a_i|\mathbf{x}_i)}{\pi_0(a_i|\mathbf{x}_i)}$$

If $c, a, \mathbf{x} \sim P(c|a, \mathbf{x})\pi_0(a|\mathbf{x})P(\mathbf{x})$

Then this can be viewed as an importance sampling formula for evaluating: $\int c P(c|a, \mathbf{x})\pi(a|\mathbf{x})P(\mathbf{x})da d\mathbf{x}$

Horvitz-Thompson Reward Estimator of a Policy

Previously we had a mapping $A(\mathbf{x})$ which returned a recommendation for user \mathbf{x} . This mapping is non-differentiable, so we introduce a “policy” $\pi(a|\mathbf{x})$. While $\pi(a|\mathbf{x})$ is mathematically a probability, but it does not represent uncertainty, it simply makes the (in this case argmax) differentiable. We then get:

$$E[U|\pi(\cdot)] = \frac{1}{N} \sum_i \frac{c_i \pi(a_i|\mathbf{x}_i)}{\pi_0(a_i|\mathbf{x}_i)}$$

If you run the epsilon greedy policy for a non-typical action $\pi_0(a|\mathbf{x}) = \frac{\epsilon}{P-1}$. This leads to massive variance..

If $c, a, \mathbf{x} \sim P(c|a, \mathbf{x})\pi_0(a|\mathbf{x})P(\mathbf{x})$

Then this can be viewed as an importance sampling formula for evaluating: $\int c P(c|a, \mathbf{x})\pi(a|\mathbf{x})P(\mathbf{x})da d\mathbf{x}$

Single Item Recommendations

Decision Rule Optimization

General Reward Estimator of a Policy

$$E[U|\pi(\cdot)] = \frac{1}{N} \sum_{i=1..N} \sum_{a'=1..P} R(a', \mathbf{x}_n) \pi(a'|\mathbf{x}_n)$$

Obviously the optimal $\pi(a|\mathbf{x})$ has the property $\pi(a^*|\mathbf{x})=1$ where $a^* = \operatorname{argmax}_a R(a, \mathbf{x})$

In practice we can obtain derivatives of $\pi(a|\mathbf{x})$ to optimize the policy/decision rule/MIPS

General Reward Estimator of a Policy

$R(a, \mathbf{x})$ can come from:

- a model
- Horvitz-Thompson
- Just do it ...

$$E[U|\pi(\cdot)] = \frac{1}{N} \sum_{i=1..N} \sum_{a'=1..P} R(a', \mathbf{x}_n) \pi(a'|\mathbf{x}_n)$$

Obviously the optimal $\pi(a|\mathbf{x})$ has the property $\pi(a^*|\mathbf{x})=1$ where $a^* = \operatorname{argmax}_a R(a, \mathbf{x})$

In practice we can obtain derivatives of $\pi(a|\mathbf{x})$ to optimize the policy/decision rule/MIPS

Optimize $\pi(a|\mathbf{x})$ when $R(a, \mathbf{x})$ is a model

If our reward model has the following parametric form:

$$R(a, \mathbf{x}) = P(c = 1|a, \mathbf{x}) = \sigma(g_{\Gamma}(\mathbf{x})^T \boldsymbol{\Psi}_a)$$

Then we already have an optimal MIPS: $f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}$, by setting
 $f_{\Xi}(\mathbf{x}) = g_{\Gamma}(\mathbf{z}), \quad \boldsymbol{\beta} = \boldsymbol{\Psi}$

So why would we optimize: $E[U|\pi(\cdot)] = \frac{1}{N} \sum_i R(a, \mathbf{x}) \pi(a|\mathbf{x})$ to obtain a MIPS if we already have one?

Reward from a model

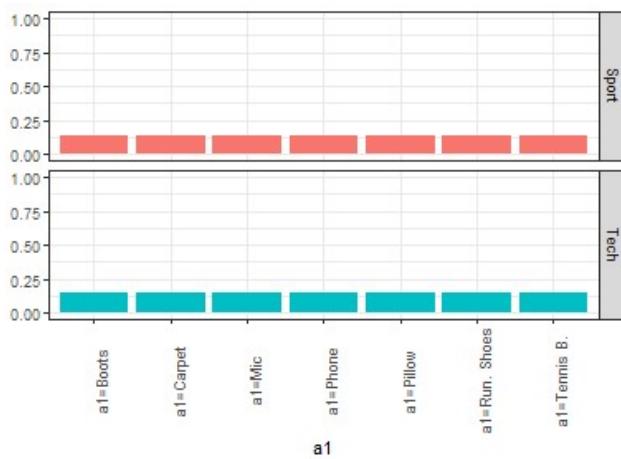
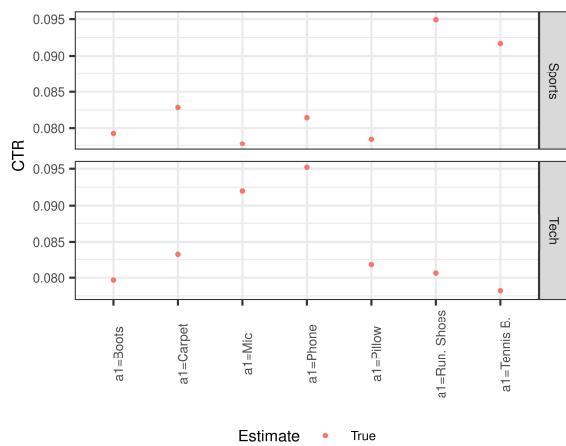
- $P(c = 1|a, \mathbf{x}) = \sigma(\mathbf{y}^T \boldsymbol{\phi} + g_\Gamma(\mathbf{z})^T \boldsymbol{\beta}_a)$ is sufficient but not necessary
- A few reasons we might want not to use:
 - $\sigma(\mathbf{y}^T \boldsymbol{\phi} + f_\Xi(\mathbf{z})^T \boldsymbol{\beta}_a)$ is a calibrated reward that correctly orders all actions.
We only need the best action
 - Imposing the above can mean the embedding dimension D for $f_\Xi(\mathbf{x})^T \boldsymbol{\beta}$ might be larger than we need
 - It might be an unreasonable engineering burden to produce perfectly optimal recommendations at all times. We might want the dimension of $\boldsymbol{\beta}$ lower than Ω
 - We might prefer a more general model $P(c = 1|a, \mathbf{x}) = \sigma(f_\Xi(\mathbf{x}, a))$

Optimizing the decision rule via a policy

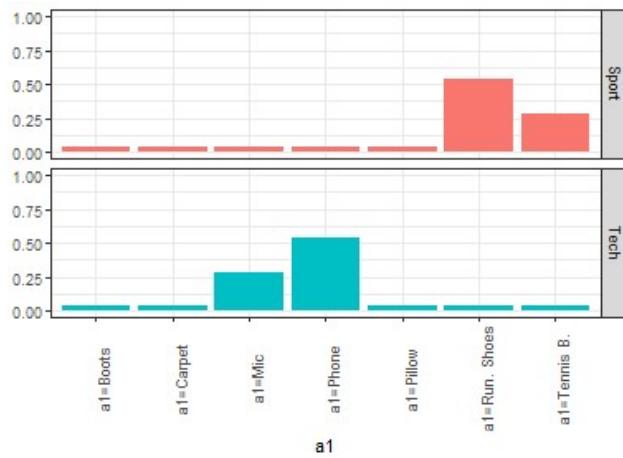
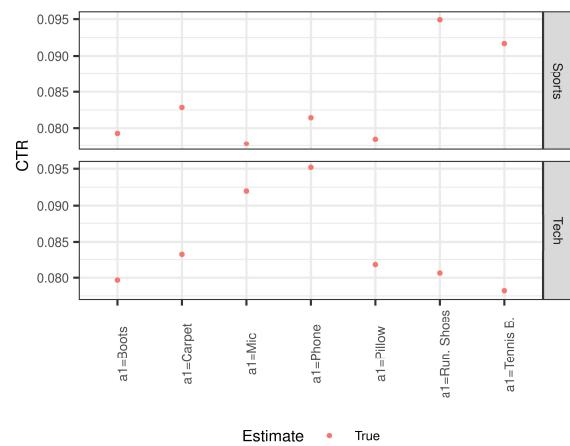
$$E[U|\pi(\cdot)] = \sum_a R(a, \mathbf{x})\pi(a|\mathbf{x})$$

If we want $\pi(a|\mathbf{x})$ to be restricted to producing a MIPS we let $\pi(a|\mathbf{x}) = \frac{\exp(f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a)}{\sum_i \exp(f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_i)}$

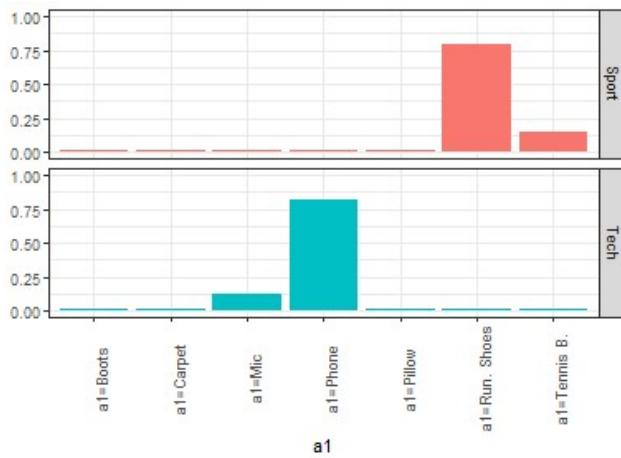
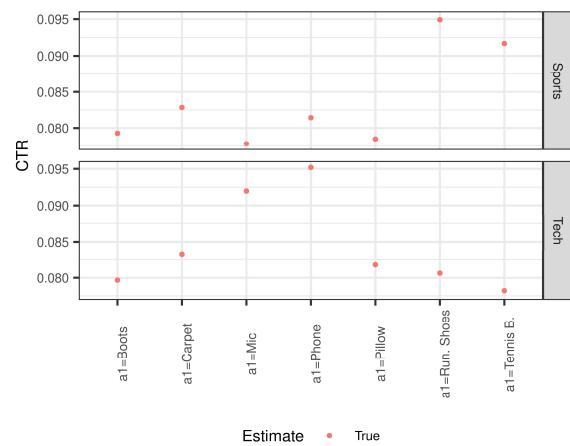
Optimizing the Policy



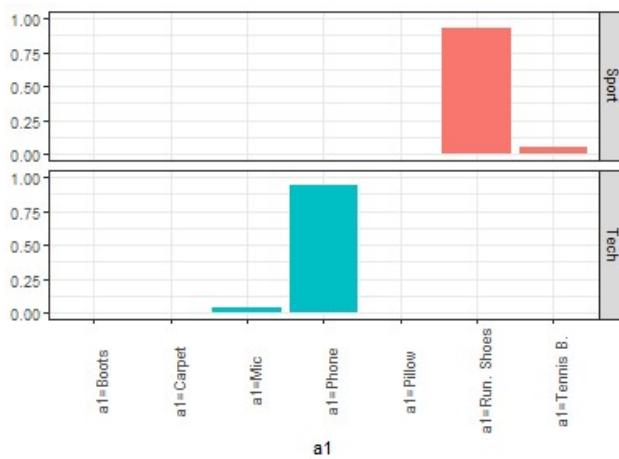
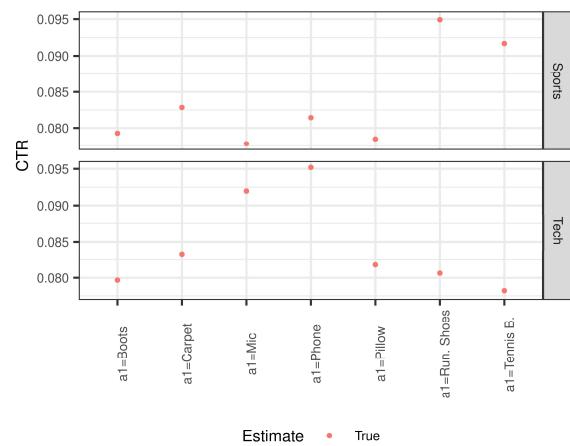
Optimizing the Policy



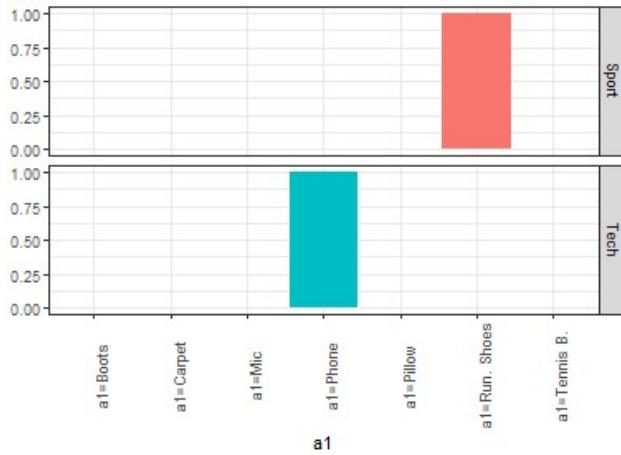
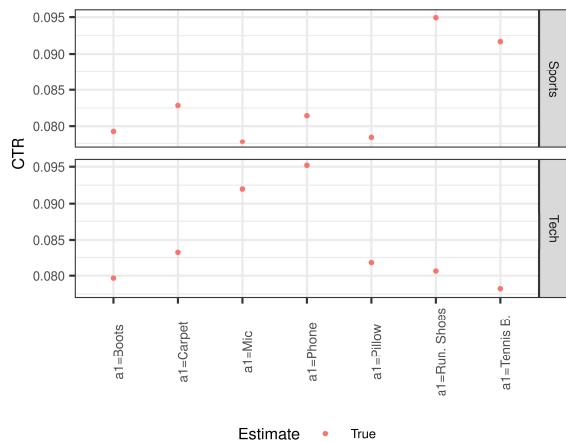
Optimizing the Policy



Optimizing the Policy

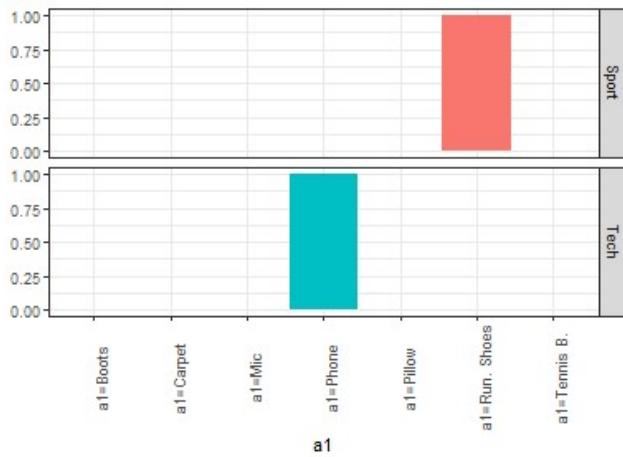
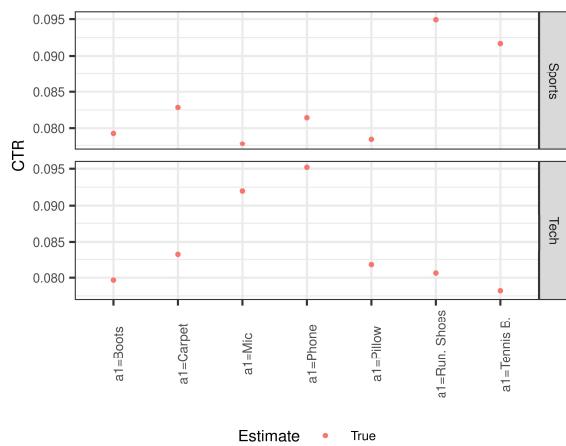


Optimizing the Policy



At convergence $\pi(a|x)$ is degenerate.

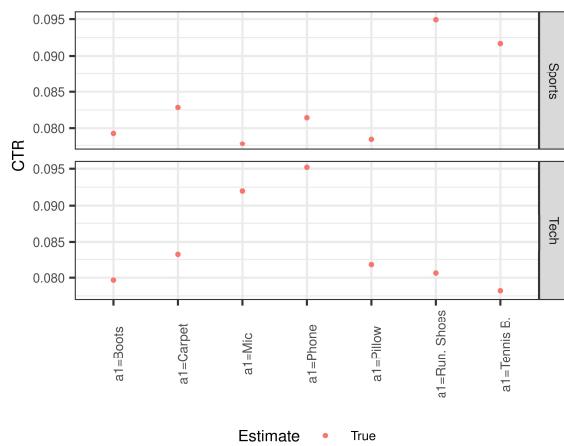
Optimizing the Policy



At convergence $\pi(a|x)$ is degenerate.

If the reward estimate is wrong, it will (possibly) select the wrong best action

Two parts to decision theoretic Reco



Reward estimation

- Model based
- Horvitz Thompson
- Just do it..



Decision rule optimization

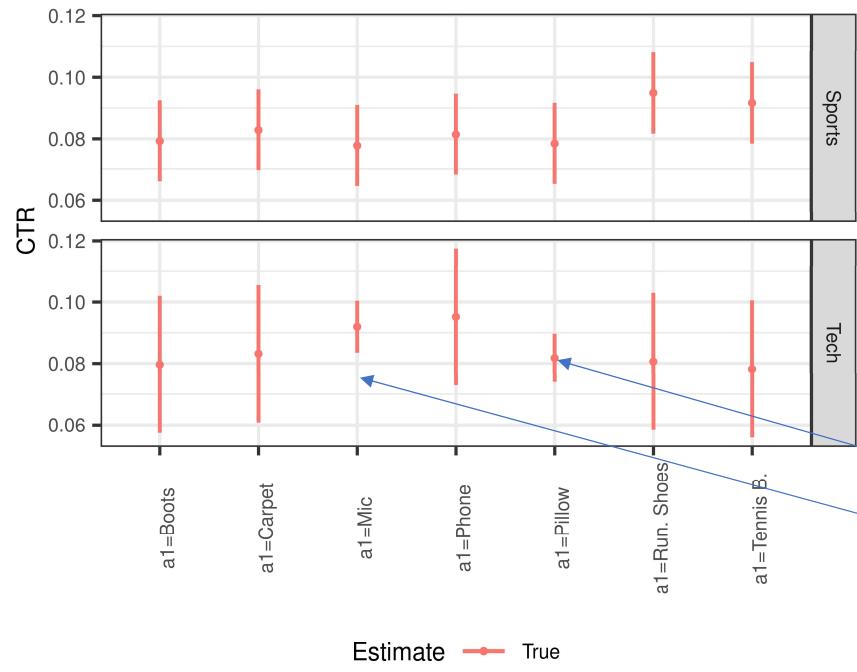
- Policy learning
- Contextual bandits
- Reinforce

Model Based Methods
Can accelerate learning
and borrow strength

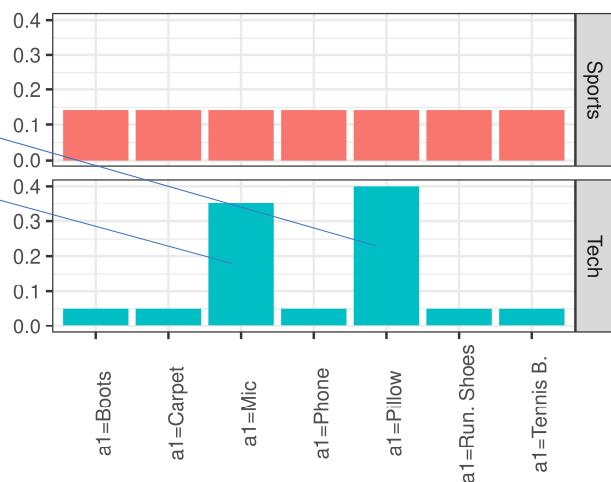
BLOB uses:

- context-context
- action-action distances
- context-action

Sakhi, O., et al. 2020



$$\pi_0(a|x)$$



Horvitz-Thompson Based Methods

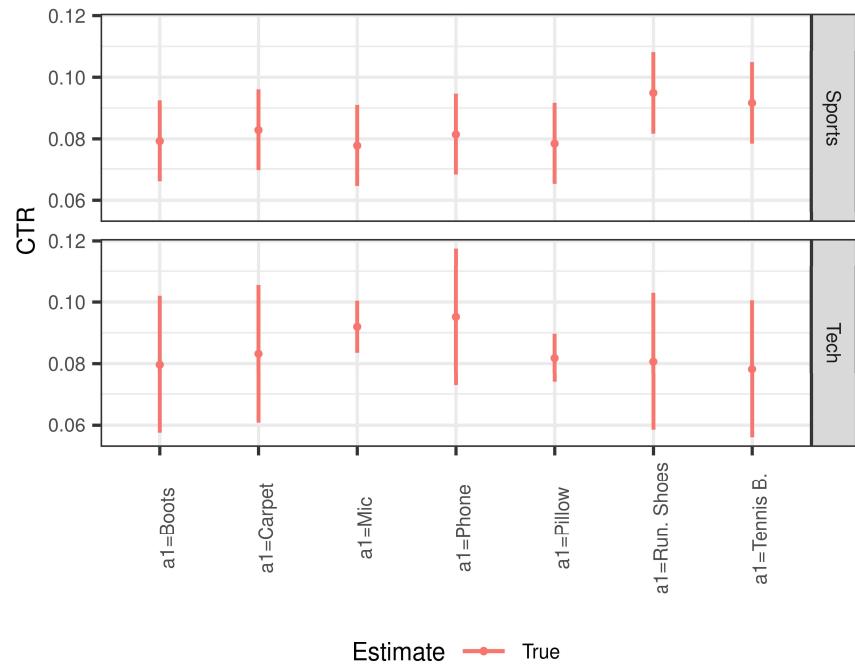
Make bias-variance trade-offs by modifying: $\frac{c_i}{\pi_0(a|x_i)}$

e.g. clipping or Stochastic Variance Penalization. “Policies Cling”

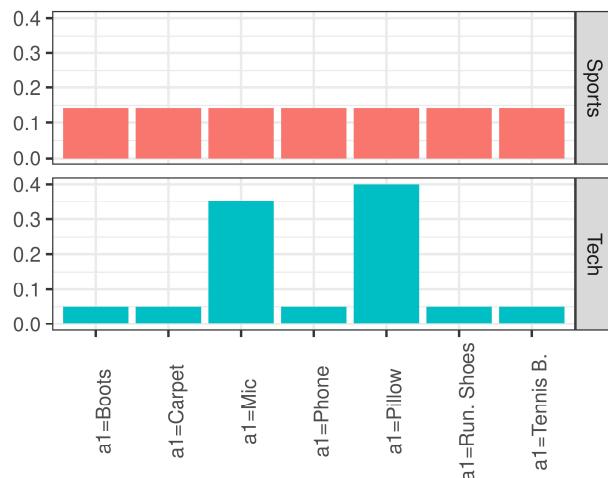
Dudík, Miroslav, et al. 2012.

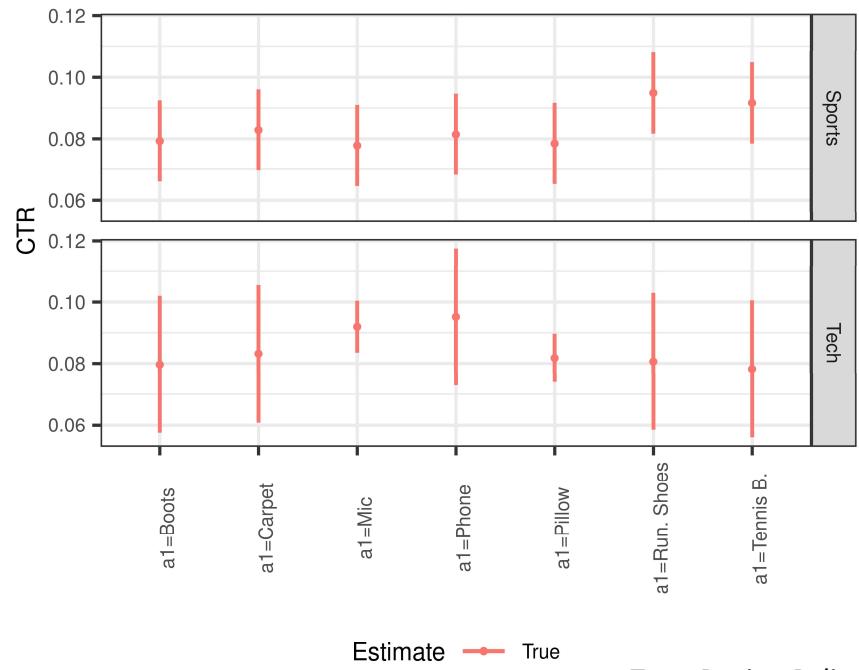
Bottou, Léon, et al. 2013

Swaminathan, Adith, and Thorsten Joachims. 2015

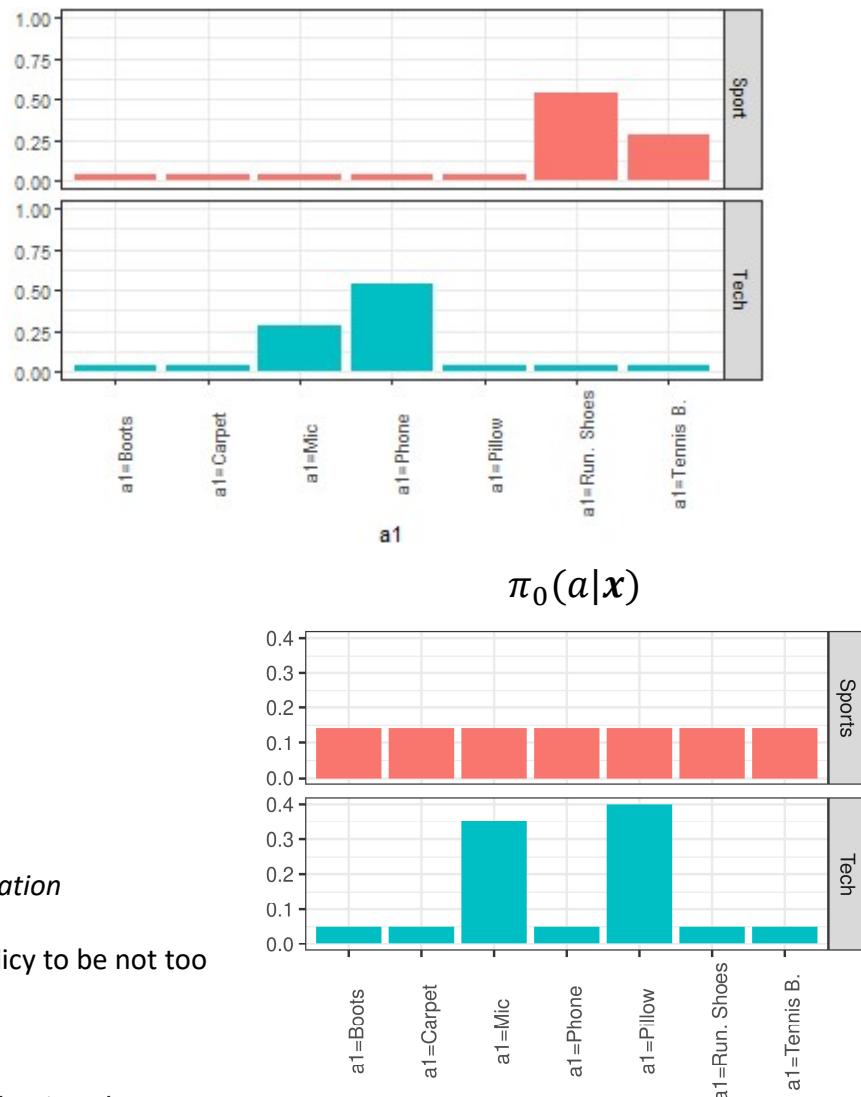


$$\pi_0(a|x)$$





Schulman, John, et al. "Trust region policy optimization."



Optimizing the MIPS with Policy Learning

$$E[U|\pi(\cdot)] = \sum_a R(a, \mathbf{z})\pi(a|\mathbf{z})$$

Works ok if the
number of actions
is very small

$$\nabla_{\Xi, \beta} E[U|\pi(\cdot)] = \nabla_{\Xi, \beta} \sum_a R(a, \mathbf{z})\pi(a|\mathbf{z})$$

Swaminathan, Adith, and Thorsten Joachims. 2015.

Joachims, Thorsten, Adith Swaminathan, and Maarten De Rijke. 2018.

Reinforce Algorithm for MIPS

$$E[R|\pi(\cdot)] = \sum_a R(a, z) \pi(a|z) = \sum_a R(a, z) \frac{\exp(f_{\Xi}(z)^T \boldsymbol{\beta}_a)}{\sum_i \exp(f_{\Xi}(z)^T \boldsymbol{\beta}_i)}$$

$$a' \sim \pi(a|z), \quad \pi(a|z) = \frac{\exp(f_{\Xi}(z)^T \boldsymbol{\beta}_a)}{\sum_i \exp(f_{\Xi}(z)^T \boldsymbol{\beta}_i)}$$

$$\nabla_{\Xi, \beta} E[R|\pi(\cdot)] = R(a', z) \nabla_{\Xi, \beta} \log \pi(a'|z)$$

On its own we don't get a speed up for using this algorithm.

From Single Action to Slates



Product
View





Product
View



Product
View





Product
View

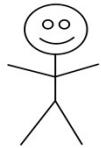


Product
View



Recommend
Slate





Product
View



Product
View



Recommend
Slate

Slate =

$$\text{argsort}(f(v_1 = \text{[image of CousCous box]}, v_2 = \text{[image of Brown Basmati Rice bag]}))^T \boldsymbol{\beta})_{1:3}$$



Product
View



Product
View



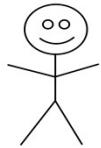
Recommend
Slate



No click

Slate =

$$\text{argsort}(f(v_1 = \text{[image of CousCous box]}, v_2 = \text{[image of CousCous box]}))^T \boldsymbol{\beta})_{1:3}$$



Product
View



Product
View



Recommend
Slate



Product
View



Slate =

$$\text{argsort}(f(v_1 = \text{[image of CousCous box]}, v_2 = \text{[image of DAAWAT rice bag]}))^T \boldsymbol{\beta})_{1:3}$$



Product
View



Product
View



Recommend
Slate



Product
View



Recommend
Slate

$$\text{Slate} = \text{argsort}(f(v_1 = \text{CousCous}, v_2 = \text{DAAWAT Brown}, v_3 = \text{DAAWAT Almond}))^T \boldsymbol{\beta}_{1:2}$$

Slate =

$$\text{argsort}(f(v_1 = \text{CousCous}, v_2 = \text{DAAWAT Brown}, v_3 = \text{DAAWAT Almond}))^T \boldsymbol{\beta}_{1:3}$$



Product
View



Product
View



Recommend
Slate



Product
View



Recommend
Slate



Slate =

$$\text{argsort}(f(v_1 = \text{CousCous}, v_2 = \text{DAAWAT}, v_3 = \text{XXX}))^T \boldsymbol{\beta}_{1:2}$$

Slate =

$$\text{argsort}(f(v_1 = \text{CousCous}, v_2 = \text{DAAWAT}, v_3 = \text{XXX}))^T \boldsymbol{\beta}_{1:3}$$



Product
View



Product
View



Recommend
Slate



Product
View



Recommend
Slate



No click

Slate =

$$\text{argsort}(f(v_1 = \text{CousCous}, v_2 = \text{DAAWAT}, v_3 = \text{XXX}))^T \boldsymbol{\beta}_{1:2}$$

Slate =

$$\text{argsort}(f(v_1 = \text{CousCous}, v_2 = \text{DAAWAT}, v_3 = \text{XXX}))^T \boldsymbol{\beta}_{1:3}$$



Product
View



Product
View



Recommend
Slate



Product
View



Recommend
Slate



Sale



Slate =

$$\text{argsort}(f(v_1 = \text{CousCous}, v_2 = \text{DAAWAT}, v_3 = \text{XXX}))^T \boldsymbol{\beta}_{1:2}$$

Slate =

$$\text{argsort}(f(v_1 = \text{CousCous}, v_2 = \text{DAAWAT}, v_3 = \text{XXX}))^T \boldsymbol{\beta}_{1:3}$$

Reward Estimation for Slates Using Horvitz-Thompson

Horvitz-Thompson Slate Reward Estimator

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \cdot c \cdot P(c | a_1, \dots, a_K, \mathbf{x}) P(\mathbf{x}) dc da_1 \dots da_K d\mathbf{x}$$

Horvitz-Thompson Slate Reward Estimator

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \cdot c \cdot P(c | a_1, \dots, a_K, \mathbf{x}) P(\mathbf{x}) dc da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \mathbb{E}[c | a_1, \dots, a_K, \mathbf{x}] P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

Horvitz-Thompson Slate Reward Estimator

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \cdot c \cdot P(c | a_1, \dots, a_K, \mathbf{x}) P(\mathbf{x}) dc da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \mathbb{E}[c | a_1, \dots, a_K, \mathbf{x}] P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \left(\sum_k \mathbb{E}[r_k | a_k \mathbf{x}] \right) P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

Horvitz-Thompson Slate Reward Estimator

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \cdot c \cdot P(c | a_1, \dots, a_K, \mathbf{x}) P(\mathbf{x}) dc da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \mathbb{E}[c | a_1, \dots, a_K, \mathbf{x}] P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \left(\sum_k \mathbb{E}[r_k | a_k \mathbf{x}] \right) P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

The fundamental assumption is this.

Li, S., Abbasi-Yadkori, Y., Kveton, B., Muthukrishnan, S., Vinay, V., & Wen, Z. 2018

Horvitz-Thompson Slate Reward Estimator

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \cdot c \cdot P(c | a_1, \dots, a_K, \mathbf{x}) P(\mathbf{x}) dc da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \mathbb{E}[c | a_1, \dots, a_K, \mathbf{x}] P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \left(\sum_k \mathbb{E}[r_k | a_k \mathbf{x}] \right) P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \left(\sum_k \pi_0(a_k | \mathbf{x}) \frac{\pi(a_k | \mathbf{x})}{\pi_0(a_k | \mathbf{x})} \mathbb{E}[r_k | a_k \mathbf{x}] \right) P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

Horvitz-Thompson Slate Reward Estimator

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \cdot c \cdot P(c | a_1, \dots, a_K, \mathbf{x}) P(\mathbf{x}) dc da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \mathbb{E}[c | a_1, \dots, a_K, \mathbf{x}] P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \left(\sum_k \mathbb{E}[r_k | a_k \mathbf{x}] \right) P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \left(\sum_k \pi_0(a_k | \mathbf{x}) \frac{\pi(a_k | \mathbf{x})}{\pi_0(a_k | \mathbf{x})} \mathbb{E}[r_k | a_k \mathbf{x}] \right) P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

Horvitz-Thompson Slate Reward Estimator

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \cdot c \cdot P(c | a_1, \dots, a_K, \mathbf{x}) P(\mathbf{x}) dc da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \mathbb{E}[c | a_1, \dots, a_K, \mathbf{x}] P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \left(\sum_k \mathbb{E}[r_k | a_k \mathbf{x}] \right) P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \left(\sum_k \pi_0(a_k | \mathbf{x}) \frac{\pi(a_k | \mathbf{x})}{\pi_0(a_k | \mathbf{x})} \mathbb{E}[r_k | a_k \mathbf{x}] \right) P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

These are marginal distributions, i.e.

$$\pi_0(a_k | \mathbf{x}) = \sum_{a'_{1'}} \dots \sum_{a'_{k-1}} \pi_0(a'_{1'}, \dots, a'_{k-1}, a_k | \mathbf{x})$$

Horvitz-Thompson Slate Reward Estimator

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \cdot c \cdot P(c | a_1, \dots, a_K, \mathbf{x}) P(\mathbf{x}) dc da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \mathbb{E}[c | a_1, \dots, a_K, \mathbf{x}] P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \pi_0(a_1, \dots, a_K | \mathbf{x}) \frac{\pi(a_1, \dots, a_K | \mathbf{x})}{\pi_0(a_1, \dots, a_K | \mathbf{x})} \left(\sum_k \mathbb{E}[r_k | a_k \mathbf{x}] \right) P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

$$E[U|\pi(\cdot)] = \int \left(\sum_k \pi_0(a_k | \mathbf{x}) \frac{\pi(a_k | \mathbf{x})}{\pi_0(a_k | \mathbf{x})} \mathbb{E}[r_k | a_k \mathbf{x}] \right) P(\mathbf{x}) da_1 \dots da_K d\mathbf{x}$$

These are marginal distributions, i.e.

$$\pi(a_k | \mathbf{x}) = \sum_{a'_1} \dots \sum_{a'_{k-1}} \pi(a'_1, \dots, a'_{k-1}, a_k | \mathbf{x})$$

Product Assumption

Independence assumptions make marginalization cheap/free... but it comes at a cost ...

$$\pi_0(a_1, \dots, a_K | \boldsymbol{x}) = \prod_k \pi_0(a_k | \boldsymbol{x})$$

$$\pi(a_1, \dots, a_K | \boldsymbol{x}) = \prod_k \pi(a_k | \boldsymbol{x})$$

Product Assumption

Independence assumptions make marginalization cheap/free... but it comes at a cost ...



$$\pi_0(a_1, \dots, a_K | x) = \prod_k \pi_0(a_k | x)$$

$$\pi(a_1, \dots, a_K | x) = \prod_k \pi(a_k | x)$$

Reward Estimation for Slates Using A Model

Probabilistic Rank and Reward

Here we outline the Probabilistic Rank and Reward Model (PRR), which has the following properties:

- We assume at most one element in the slate is clicked
 - This differs with IIPS, Position Based Model, Cascade Model etc
- PRR includes user engagement features and recommendation features
- PRR incorporates:
 - was the slate clicked? *bandit approach*
 - which item was clicked? *IIPS, position-based model*
- Gives a direct MIPS

Details in: [Aouali et al. 2022](#)

Preference Learning vs Reward

- **Ranking model** (like BPR)
- Trained on clicked-banners only
- Not reward optimizing (it doesn't see it)



Reward model

- Trained on all banners
- Returns calibrated CTR



Preference Learning vs Reward

- **Ranking model** (like BPR)
- Trained on clicked-banners only
- Not reward optimizing (it doesn't see it)



Reward model

- Trained on all banners
- Returns calibrated CTR



PRR combines both!

The Probabilistic Rank and Reward model (our approach)

$$\bar{c}, r_1, \dots, r_K | a_1, \dots, a_K, \mathbf{y}, \mathbf{z} \sim \text{categorical} \left(\frac{\theta_0}{Z}, \frac{\theta_1}{Z}, \dots, \frac{\theta_K}{Z} \right), \quad Z = \sum \theta_k$$

\uparrow

$$\theta_0 = \exp(\mathbf{y}^T \boldsymbol{\phi})$$

\mathbf{y} – User engagement features (or bidding features)

$\boldsymbol{\phi}$ – nuisance parameters

The Probabilistic Rank and Reward model (our approach)

$$\bar{c}, r_1, \dots, r_K | a_1, \dots, a_K, \mathbf{y}, \mathbf{z} \sim \text{categorical} \left(\underbrace{\frac{\theta_0}{Z}, \frac{\theta_1}{Z}, \dots, \frac{\theta_K}{Z}}_{\gamma} \right), \quad Z = \sum \theta_k$$

$$\theta_k = \exp(g_\Gamma(\mathbf{z})^T \boldsymbol{\Psi}_{a_k}) \exp(\gamma_k) + \exp(\alpha_k)$$

\mathbf{z} – context

$g_\Gamma(\mathbf{z})$ – contextual embedding

$\boldsymbol{\Psi}_a$ – Item embedding for item a

γ, α – multiplicative and additional position bias

Reminder

$a_1, \dots, a_K = \text{argsort}(f_\Xi(\mathbf{z})^T \boldsymbol{\beta})$

$f_\Xi = g_\Gamma$ and $\boldsymbol{\beta} = \boldsymbol{\Psi}$

Reward and Rank only variants

Reward model

$$\bar{c}, c | a_1, \dots, a_K, y, z \sim \text{categorical} \left(\frac{\theta_0}{Z}, \sum \frac{\theta_k}{Z} \right), \quad Z = \sum_{k=0}^K \theta_k$$

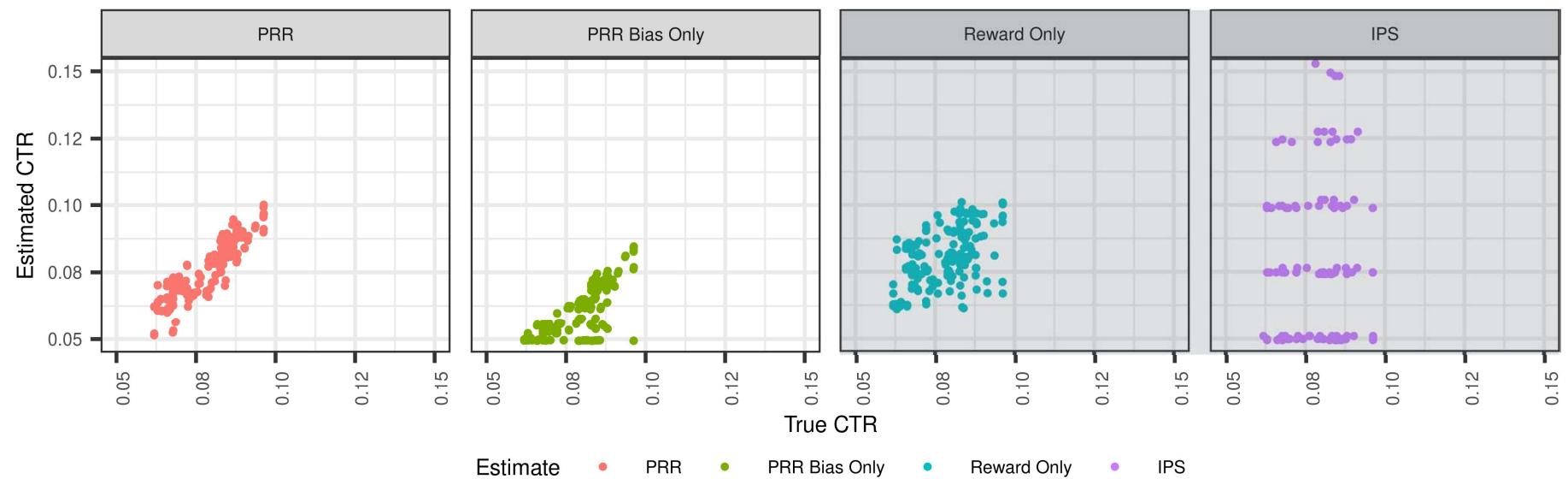
Ranking model

$$r_1, \dots, r_K | c = 1, a_1, \dots, a_K, y, z \sim \text{categorical} \left(\frac{\theta_1}{Z}, \dots, \frac{\theta_K}{Z} \right), \quad Z = \sum_{k=1}^K \theta_k$$

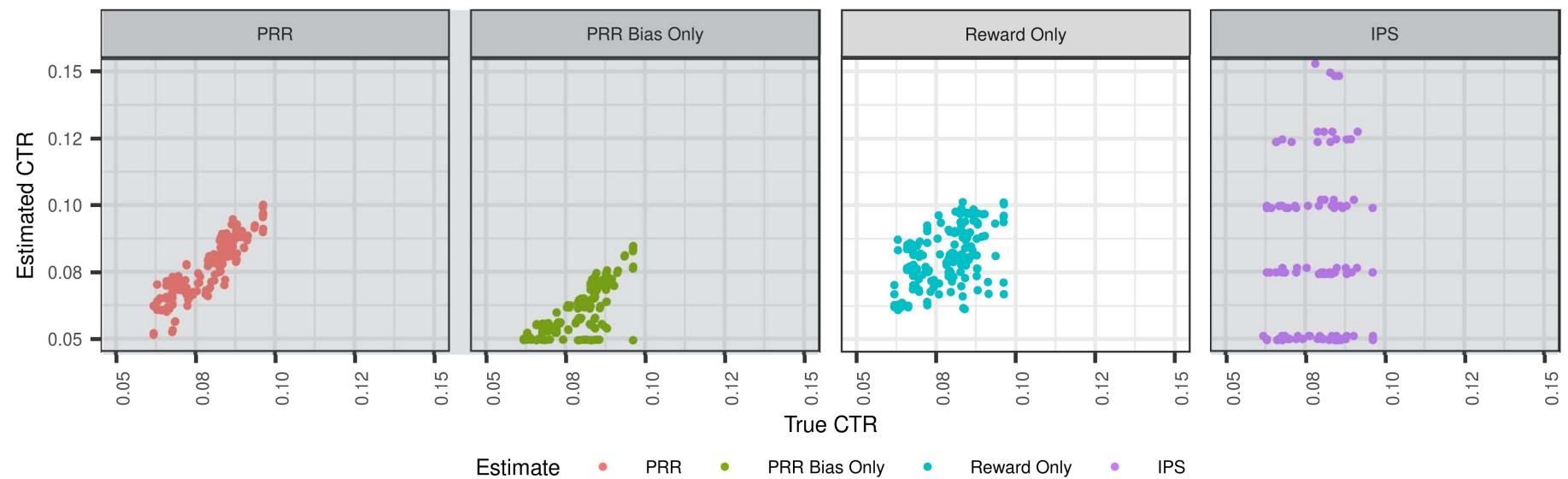
$$\theta_0 = \exp(\mathbf{y}^T \boldsymbol{\phi})$$

$$\theta_K = \exp(g_\Gamma(\mathbf{z})^T \boldsymbol{\Psi}_{ak}) \exp(\gamma_k) + \exp(\alpha_k)$$

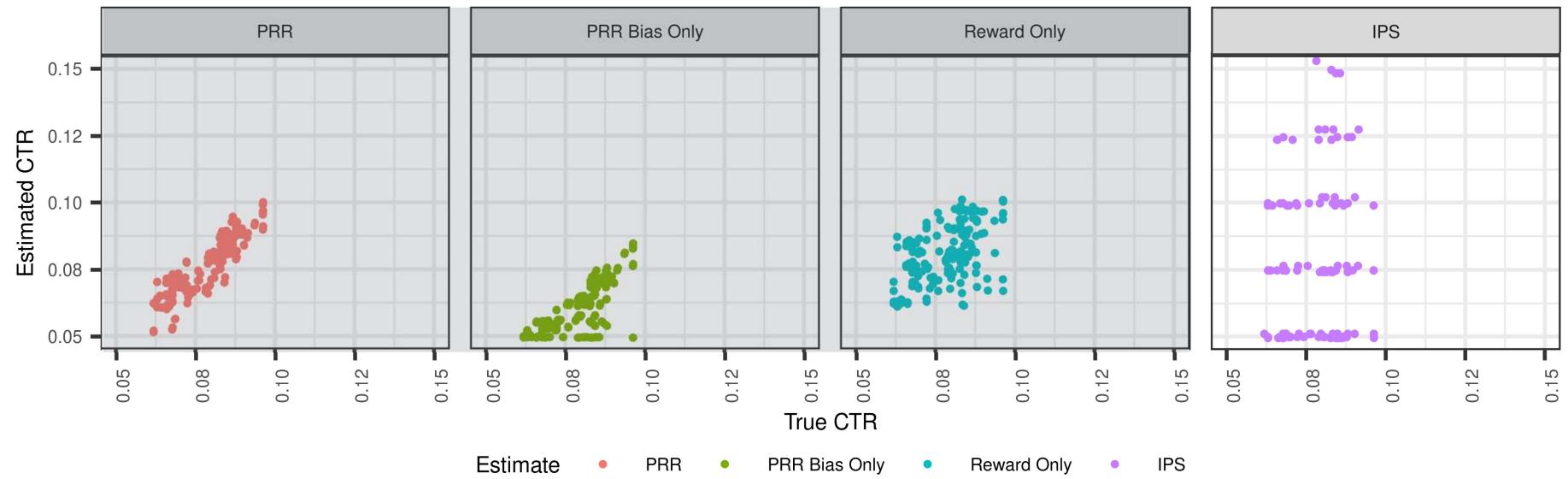
Calibration on toy example



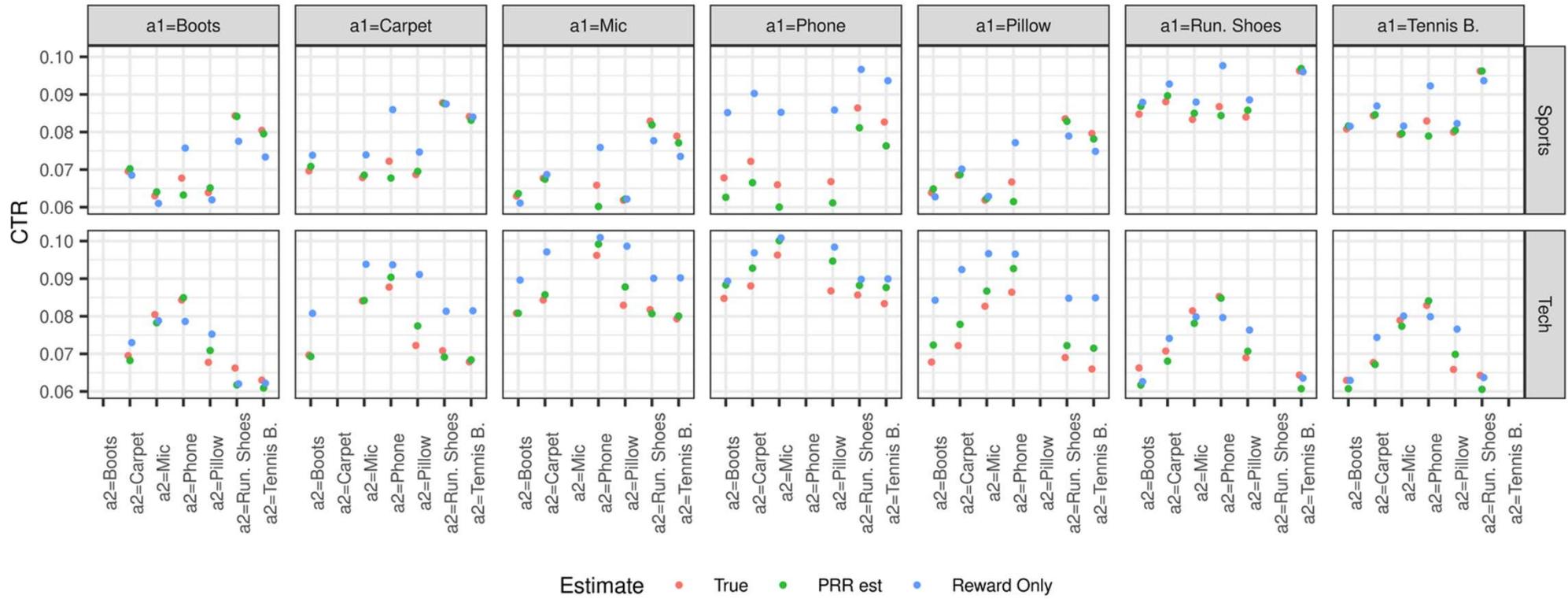
Calibration on toy example



Calibration on toy example



Click-Through Rate Estimates



Slate Recommendations

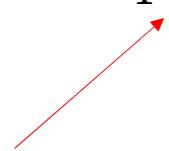
Decision Rule Optimization

Reinforce Algorithm for Slates

$$\mathbb{E}[U|\pi(\cdot)] = \sum_{a_1} \dots \sum_{a_K} R(a_1, \dots, a_K, \mathbf{z}) \pi(a_1, \dots, a_K | \mathbf{z})$$

$$\nabla_{\Xi, \beta} \mathbb{E}[U|\pi(\cdot)] = \nabla_{\Xi, \beta} \sum_{a_1} \dots \sum_{a_K} R(a_1, \dots, a_K, \mathbf{z}) \pi(a_1, \dots, a_K | \mathbf{z})$$

This is now totally
infeasible



Reinforce Algorithm for Slates

$$\mathbb{E}[U|\pi(\cdot)] = \sum_a R(a_1, \dots, a_K, \mathbf{x})\pi(a_1, \dots, a_K | \mathbf{x})$$

$$a'_1, \dots, a'_K \sim \pi(\cdot | \mathbf{x})$$

$$\nabla_{\Xi, \beta} \mathbb{E}[U|\pi(\cdot)] = R(a'_1, \dots, a'_K, \mathbf{x}) \nabla_{\Xi, \beta} \log \pi(a'_1, \dots, a'_K | \mathbf{x})$$

If we can sample and evaluate $\pi(\cdot | \mathbf{x})$ then we can proceed...

Independence Assumption



$$\pi(a_1, \dots, a_K | x) = \prod_k \pi(a_k | x)$$

What if the best two recommendations are rice and beer

Firstly, we don't want:

$$\pi(a = \text{rice} | x) = 1$$

as this forces duplication for large slates i.e.

$$\pi(a_1 = \text{rice}, a_2 = \text{rice} | x) = 1$$

$$f_{\Xi}(\mathbf{z})^T \boldsymbol{\beta}_{\text{rice}} \gg f_{\Xi}(\mathbf{z})^T \boldsymbol{\beta}_i, \forall i \neq \text{rice}$$

Independence Assumption



$$\pi(a_1, \dots, a_K | x) = \prod_k \pi(a_k | x)$$

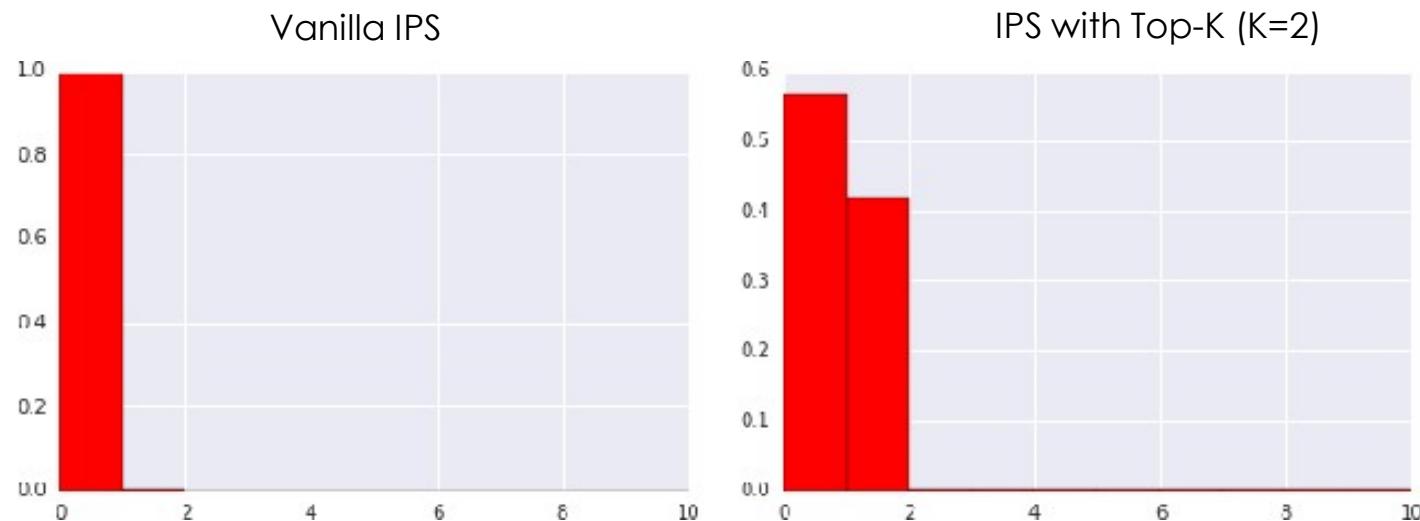
What if the best two recommendations are rice and beer

Firstly, we don't want:

$$\pi(a = \text{rice} | x) = 1$$

Chen, 2019 proposes a solution to this which keeps the probability spread out over the Top-K items.

Independence Assumption With Top-K



From Top-K Off-Policy Correction for a REINFORCE Recommender System [1]

The Top-K correction
modifies the loss to stop
 $\pi(a_k|x)$ converging to a
degenerate distribution

$$\pi(a_1, \dots, a_K | \mathbf{x}) = \prod \pi(a_k | \mathbf{x})$$

Chen, et al. 2019

Independence Assumption



$$\pi(a_1, \dots, a_K | x) = \prod_k \pi(a_k | x)$$

What if the best two recommendations are rice and beer

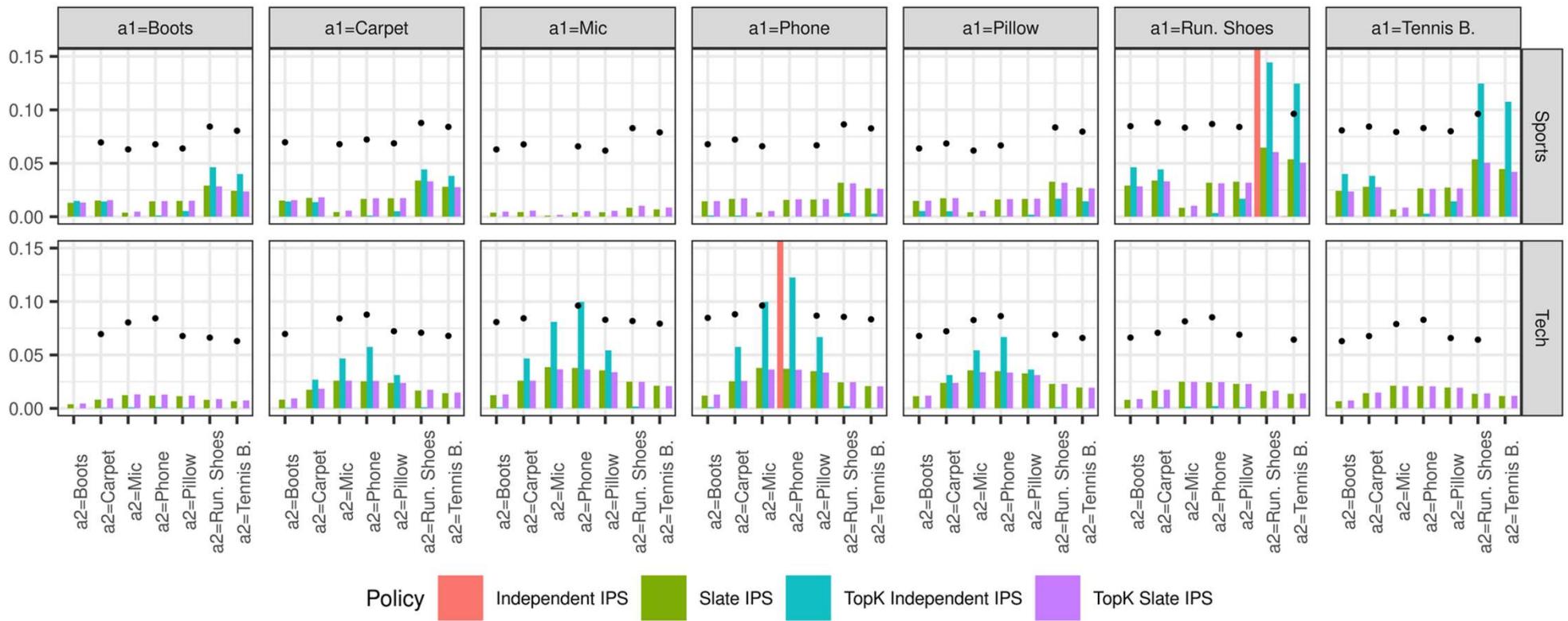
$$\pi(a = \text{rice} | x) = 0.51$$

$$\pi(a = \text{beer} | x) = 0.49$$

Then

$$\begin{aligned}\pi(a_1 = \text{rice}, a_2 = \text{rice} | x) &= 0.2601 \\ \pi(a_1 = \text{rice}, a_2 = \text{beer} | x) &= 0.2499 \\ \pi(a_1 = \text{beer}, a_2 = \text{rice} | x) &= 0.2499 \\ \pi(a_1 = \text{beer}, a_2 = \text{beer} | x) &= 0.2401\end{aligned}$$

Policy Learning



Placket-Luce Policy Formulation

Sampling without replacement is called Placket-Luce

$$\pi(a_1, \dots, a_K | \mathbf{z}) = \frac{e^{f_{\Xi}(\mathbf{z})^T \boldsymbol{\beta}_{a_1}}}{\sum e^{f_{\Xi}(\mathbf{z})^T \boldsymbol{\beta}_i}} \cdot \frac{e^{f_{\Xi}(\mathbf{z})^T \boldsymbol{\beta}_{a_2}}}{e^{f_{\Xi}(\mathbf{z})^T \boldsymbol{\beta}_{a_1}} - \sum e^{f_{\Xi}(\mathbf{z})^T \boldsymbol{\beta}_i}} \cdots \frac{e^{f_{\Xi}(\mathbf{z})^T \boldsymbol{\beta}_{a_K}}}{e^{f_{\Xi}(\mathbf{z})^T \boldsymbol{\beta}_{a_1}} + \dots + e^{f_{\Xi}(\mathbf{z})^T \boldsymbol{\beta}_{a_{K-1}}} - \sum e^{f_{\Xi}(\mathbf{z})^T \boldsymbol{\beta}_i}}$$

Reinforce Algorithm for MIPS

$$E[U|\pi(\cdot)] = \sum_a R(a_1, \dots, a_K, \mathbf{x})\pi(a_1, \dots, a_K | \mathbf{x})$$

$$a'_1, \dots, a'_K \sim \pi(\cdot | \mathbf{x})$$

$$\nabla_{\Xi, \beta} E[U|\pi(\cdot)] = R(a'_1, \dots, a'_K, \mathbf{x}) \log \nabla_{\Xi, \beta} \pi(a'_1, \dots, a'_K | \mathbf{x})$$

Simulated A/B Tests – Slate Sizes

Algorithm	K=2	K=6	K=18
Oracle	$72.62 \pm 0.45\%$	$64.43 \pm 0.48\%$	$27.08 \pm 0.44\%$
Random	$3.39 \pm 0.18\%$	$7.72 \pm 0.27\%$	$5.24 \pm 0.22\%$
PRR	$50.47 \pm 0.5\%$	$61.70 \pm 0.49\%$	$26.7 \pm 0.44\%$
PRR (Bias Only)	$36.0 \pm 0.48\%$	$60.12 \pm 0.49\%$	$25.44 \pm 0.44\%$
Ranking Model	$35.79 \pm 0.48\%$	$48.88 \pm 0.5\%$	$26.31 \pm 0.44\%$
Reward Model	$42.04 \pm 0.49\%$	$22.28 \pm 0.42\%$	$5.43 \pm 0.23\%$
Slate IPS	$20.02 \pm 0.4\%$	$26.67 \pm 0.44\%$	$3.31 \pm 0.18\%$
Top-K Slate IPS	$20.02 \pm 0.4\%$	$17.84 \pm 0.38\%$	$3.31 \pm 0.18\%$
IIPS	$40.42 \pm 0.49\%$	$48.91 \pm 0.5\%$	$25.36 \pm 0.44\%$
Top-K IIPS	$40.86 \pm 0.49\%$	$52.77 \pm 0.5\%$	$25.6 \pm 0.44\%$

Table 3: Simulated A/B test with varying slate size. $P = 1k$, $N_{\text{Train}} = 10k$

Simulated A/B Tests – Scaling

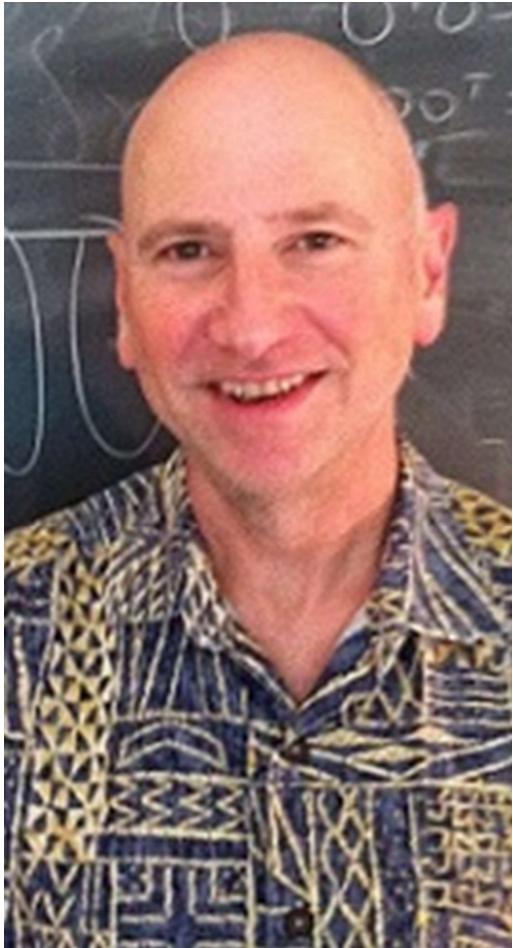
Algorithm	P=1k	P=5k	P=10k
Oracle	$49.21 \pm 0.50\%$	$76.94 \pm 0.42\%$	$78.01 \pm 0.41\%$
Random	$3.65 \pm 0.19\%$	$7.0 \pm 0.26\%$	$5.68 \pm 0.23\%$
PRR	$42.63 \pm 0.49\%$	$58.01 \pm 0.49\%$	$57.18 \pm 0.49\%$
PRR (Bias only)	$42.88 \pm 0.49\%$	$49.98 \pm 0.5\%$	$44.31 \pm 0.5\%$
Ranking Model	$26.15 \pm 0.44\%$	$50.57 \pm 0.5\%$	$49.72 \pm 0.5\%$
Reward Model	$2.83 \pm 0.17\%$	$38.35 \pm 0.49\%$	$19.14 \pm 0.39\%$
Slate IPS	$20.93 \pm 0.41\%$	$29.19 \pm 0.45\%$	$18.01 \pm 0.38\%$
Top-K Slate IPS	$27.23 \pm 0.45\%$	$25.38 \pm 0.44\%$	$15.76 \pm 0.36\%$
IIPS	$36.56 \pm 0.48\%$	$51.47 \pm 0.5\%$	$34.85 \pm 0.48\%$
Top-K IIPS	$43.01 \pm 0.5\%$	$53.96 \pm 0.5\%$	$37.24 \pm 0.48\%$

Table 4: Simulated A/B test with varying catalog size. $K = 5$, $N_{\text{Train}} = 10k$

Simulated A/B Tests – 1M Products

Algorithm	CTR	Time/epoch (s)
Oracle	$94.84 \pm 0.31\%$	-
Random	$6.25 \pm 0.34\%$	-
PRR	$67.19 \pm 0.66\%$	4.36
PRR (Bias only)	$46.87 \pm 0.71\%$	4.36
Ranking Model	$52.13 \pm 0.71\%$	2.81
Reward Model	$31.37 \pm 0.66\%$	4.37
Slate IPS	$3.11 \pm 0.25\%$	68.38
Top-K SlateIPS	$3.11 \pm 0.25\%$	69.81
IIPS	$42.22 \pm 0.7\%$	68.60
Top-K IIPS	$48.58 \pm 0.71\%$	69.56

Table 6: Simulated A/B test for 1M products. $N_{\text{Train}} = 10k$, $K = 5$, batch size 516. Note that Ranking and IPS methods are only using 6% of the dataset.



A Great Controversy in Statistics

Christopher Sims and Larry Wasserman had an “epic debate” on the Horvitz-Thompson estimator

Around if it causes a problem for likelihood/Bayesian inference

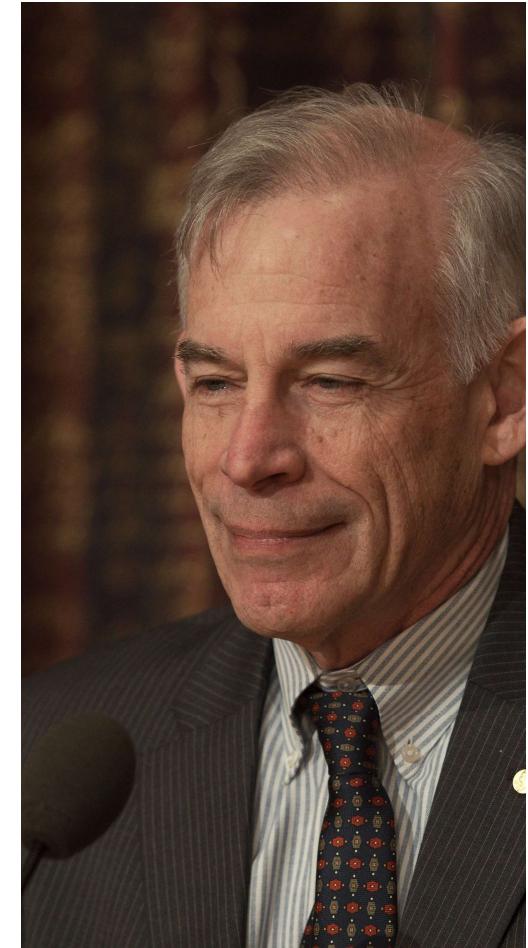
Bayesians/Likelihoodists should not use the propensity score..

PRR + BLOB ignore the propensity score and outperform other methods

Debate and discussion:

[Bayesian Causal Inference for Real World Interactive Systems
\(bcirwis2021.github.io\)](https://bcirwis2021.github.io)

[Christopher Sims - Large Parameter Spaces and Weighted Data: A Bayesian Perspective - YouTube](#)



Notebook 1 (return to)

[https://colab.research.google.com/github/otmhi/Reward-Optimizing-Reco/blob/main/Reward Optimizing Slate Recommendation with DL and MIPS Part 1.ipynb](https://colab.research.google.com/github/otmhi/Reward-Optimizing-Reco/blob/main/Reward%20Optimizing%20Slate%20Recommendation%20with%20DL%20and%20MIPS%20Part%201.ipynb)

“Technology”: SNIPS Reinforce

Optimizing the Decision Rule Quickly

Optimizing the objective

Our objective is:

$$U_{\Xi,\beta}(\mathbf{x}) = \mathbb{E}_{a \sim \pi_{\Xi,\beta}(\cdot|\mathbf{x})}[R(a, \mathbf{x})]$$

$$\nabla_{\Xi,\beta} U_{\Xi,\beta}(\mathbf{x}) = \mathbb{E}_{a \sim \pi_{\Xi,\beta}(\cdot|\mathbf{x})}[R(a, \mathbf{x}) \nabla_{\Xi,\beta} \log \pi_{\Xi,\beta}(a|\mathbf{x})]$$

There are mainly two problems with this gradient :

- $\log \pi_{\theta}(a|x)$ involves the computation of the normalizing constant $\rightarrow O(P)$
- If we want to compute the exact expectation $\rightarrow O(P)$ (Infeasible)
- If we want to approximate the gradient by sampling $\rightarrow O(P)$

Everything scales **at least linearly in P** which can be very costly in large catalog sizes.

Policy Optimization is not Maximum Likelihood

$$\log P(a|x, \beta, \Xi) = f_{\Xi}(x)^T \beta_a - \log \sum_{a'} \exp(f_{\Xi}(x)^T \beta_{a'})$$

$$\nabla_{\beta, \Xi} \log P(a|x, \beta, \Xi) = \nabla_{\beta, \Xi} f_{\Xi}(x)^T \beta_a - \mathbb{E}_P(a'|x, \beta, \Xi) [\nabla_{\beta, \Xi} f_{\Xi}(x)^T \beta_{a'}]$$

Lots of work on how to solve this problem! e.g.

Bengio, Yoshua, and Jean-Sébastien Senécal. 2003.

Ruiz, Francisco, et al. 2018.

Rawat, Ankit Singh, et al. 2019.

Optimizing the objective

Our objective is:

$$U_{\Xi,\beta}(\mathbf{x}) = \mathbb{E}_{A \sim \pi_{\Xi,\beta}(\cdot|\mathbf{x})}[R(a, \mathbf{x})]$$

$$\nabla_{\Xi,\beta} U_{\Xi,\beta}(\mathbf{x}) = \mathbb{E}_{a \sim \pi_{\Xi,\beta}(\cdot|\mathbf{x})}[R(a, \mathbf{x}) \nabla_{\Xi,\beta} \log \pi_{\Xi,\beta}(a|\mathbf{x})]$$

$$\nabla_{\Xi,\beta} U_{\Xi,\beta}(\mathbf{x}) = Cov_{a \sim \pi_{\theta}(\cdot|\mathbf{x})}[R(a, \mathbf{x}), \nabla_{\Xi,\beta} \exp f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a]$$

With $Cov[Y, \mathbf{Z}] = \mathbb{E}[(Y - \mathbb{E}[Y])(\mathbf{Z} - \mathbb{E}[\mathbf{Z}])]$

Proposal Selection

Additional Assumption

$$\pi(a|x) = \frac{\exp(f_{\Xi}(x)^T \beta)}{\sum_{a'} \exp(f_{\Xi}(x)^T \beta_{a'})}$$

Assume β is fixed and we only optimize $f_{\Xi}(x)^T$. This structure allows fast search/argsort $O(\log P)$.

MIPS give us $\alpha_S(x) = \text{argsort}(f_{\Xi}(x)^T \beta)_{[1, \dots, S]}$ the top (approximate) items of the policy

We define :

$$Q_{\epsilon, K}(a|x) = \begin{cases} \frac{\epsilon}{P} + (1 - \epsilon)\kappa(a|x) & \text{if } a \in \alpha_K(x) \\ \frac{\epsilon}{P} & \text{else} \end{cases} \quad \text{with } \kappa_{\theta}(a|x) = \frac{\exp(f_{\Xi}(x)^T \beta_a)}{\sum_{a' \in \alpha_K(x)} \exp(f_{\Xi}(x)^T \beta_{a'})}$$

Self Normalized Importance Sampling

$$\tilde{\pi}(a|x) = \exp(f_{\Xi}(x)^T \boldsymbol{\beta}_a)$$

$$Q_{\epsilon,K}(a|X) = \begin{cases} \frac{\epsilon}{P} + (1 - \epsilon)\kappa(a|x) & \text{if } a \in \alpha_K(x) \\ \frac{\epsilon}{P} & \text{else} \end{cases} \quad \text{with } \kappa_{\Xi,\beta}(a|x) = \frac{\exp(f_{\Xi}(x)^T \boldsymbol{\beta}_a)}{\sum_{a' \in \alpha_K(x)} \exp(f_{\Xi}(x)^T \boldsymbol{\beta}_{a'})}$$

$$\begin{aligned} \tilde{w}_i &= \frac{\tilde{\pi}(A_i|x)}{Q(A_i|x)}, & w_i &= \frac{\tilde{w}_i}{\sum \tilde{w}_j} \\ \bar{Y} &= \sum_{i=1}^N w_i Y_i, & \bar{Z} &= \sum_{i=1}^N w_i Z_i \end{aligned}$$

$$\begin{aligned} \nabla_{\Xi,\beta} U_{\Xi,\beta}(x) &= Cov_{a \sim \pi_{\Xi,\beta}(\cdot|x)}[R(a,x), \nabla_{\Xi,\beta} \exp f_{\Xi}(x)^T \boldsymbol{\beta}_a], \\ \text{with } Cov[Y, Z] &= \mathbb{E}[(Y - \mathbb{E}[Y])(Z - \mathbb{E}[Z])] \end{aligned}$$

$$\text{return } \widehat{SNIS}_N = \sum_{i=1}^N w_i (Y_i - \bar{Y})(Z_i - \bar{Z})$$

Experiments

We test the method on a session completion task : X observed session, Y complementary items.

This can be cast into an offline bandit framework where our policy π_θ takes the observed part X as a context, recommends an item a and receives the binary reward $\hat{r}(a, X) = 1[a \in Y]$.

We use X to build item embeddings β of dimension $L \ll P$.

Define x_{emb} as the item embeddings average in session X.

$$\pi(a|x) = \frac{\exp(f_\Xi(x)^T \beta_a)}{\sum_{a'} \exp(f_\Xi(x)^T \beta_{a'})}$$

	Catalog size	Number of users
Twitch	750K	500K
GoodReads	1.23M	300K

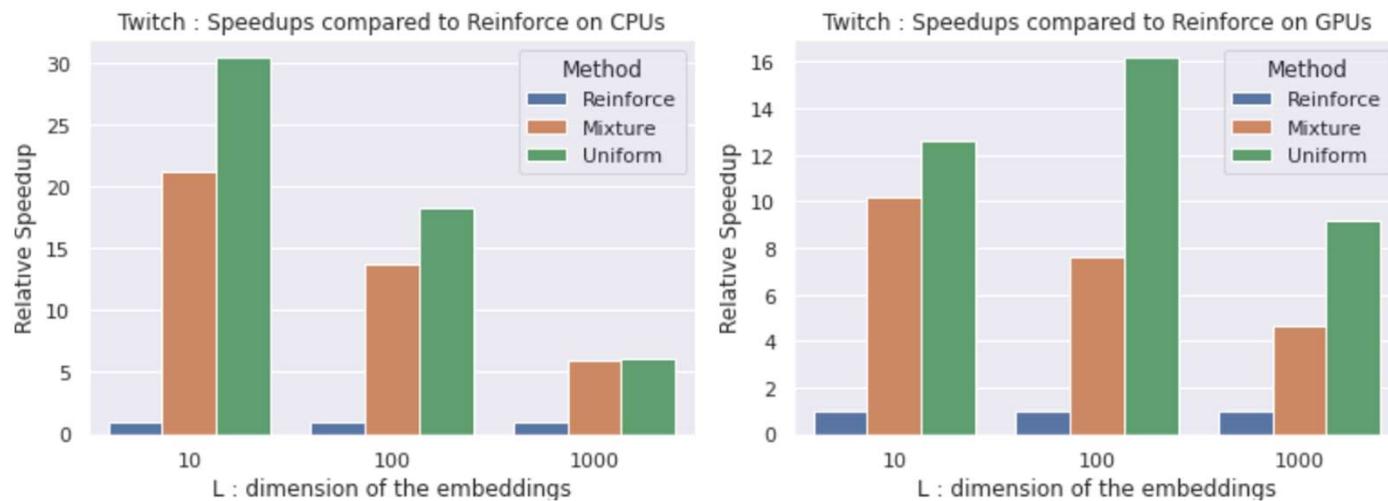
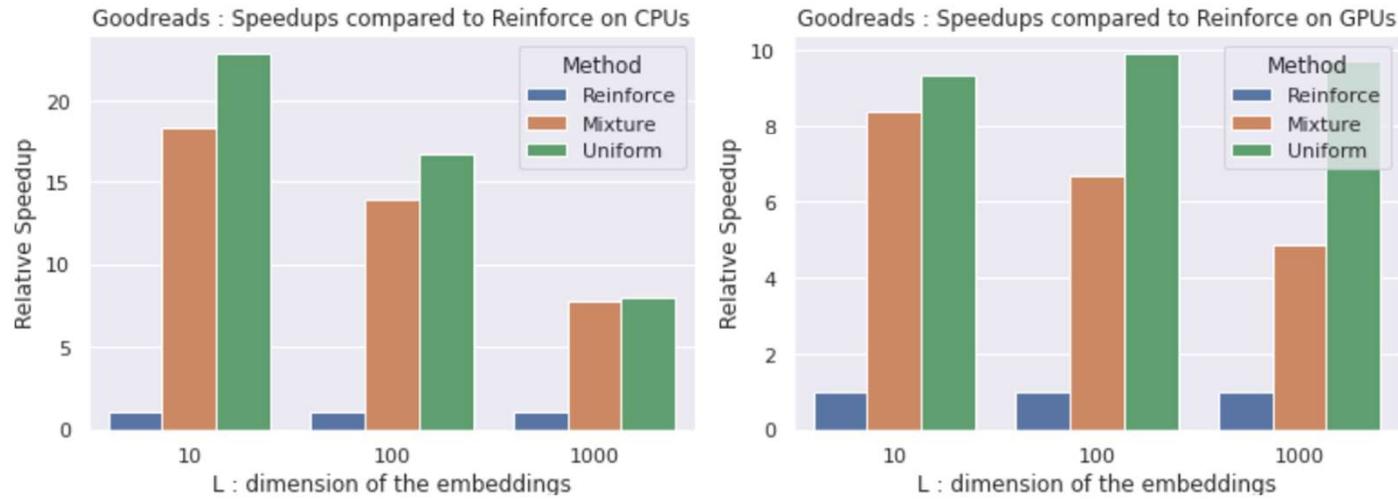


Fig. 1. The Twitch Dataset : The speedup of the proposed algorithms on both CPU and GPU devices



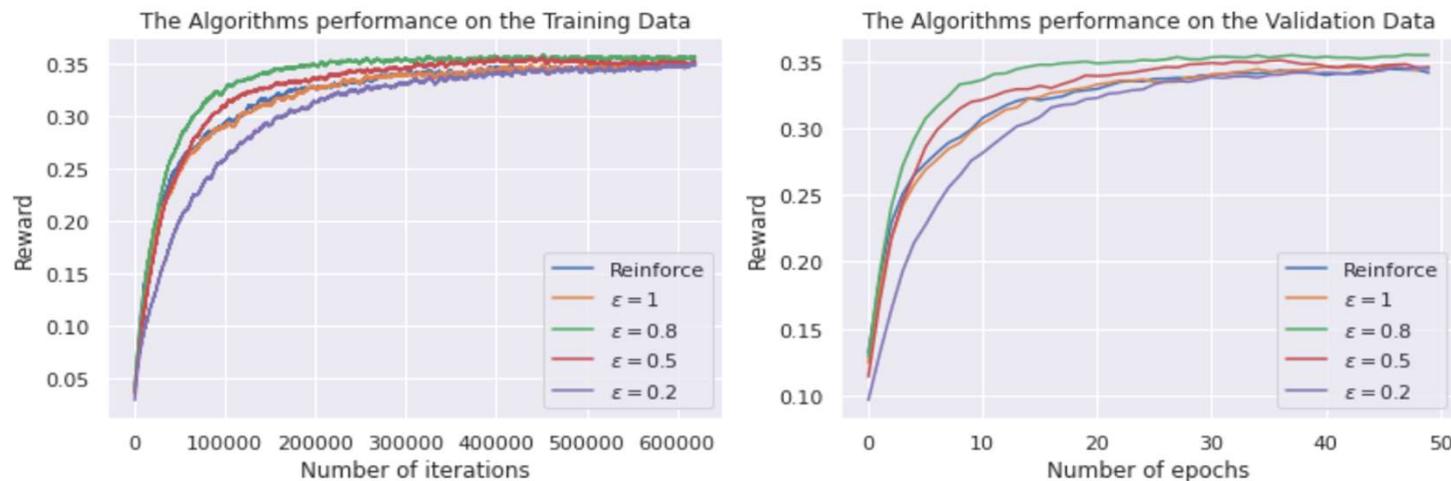


Fig. 3. The Twitch Dataset : ϵ effect on the policy training and validation

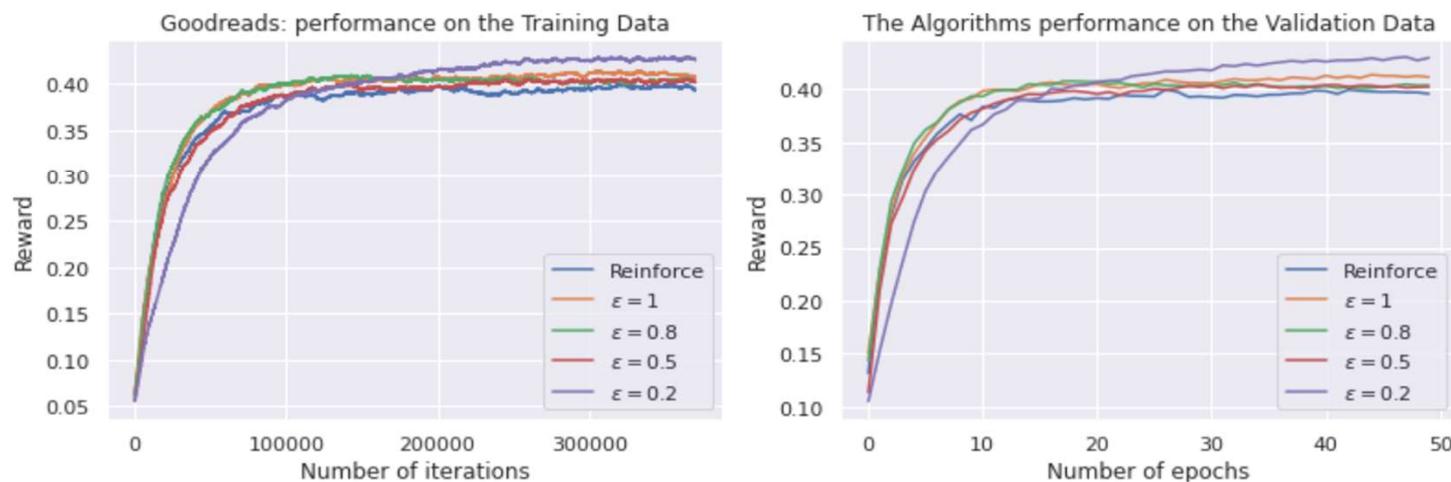


Fig. 4. The Goodreads Dataset : ϵ effect on the policy training and validation

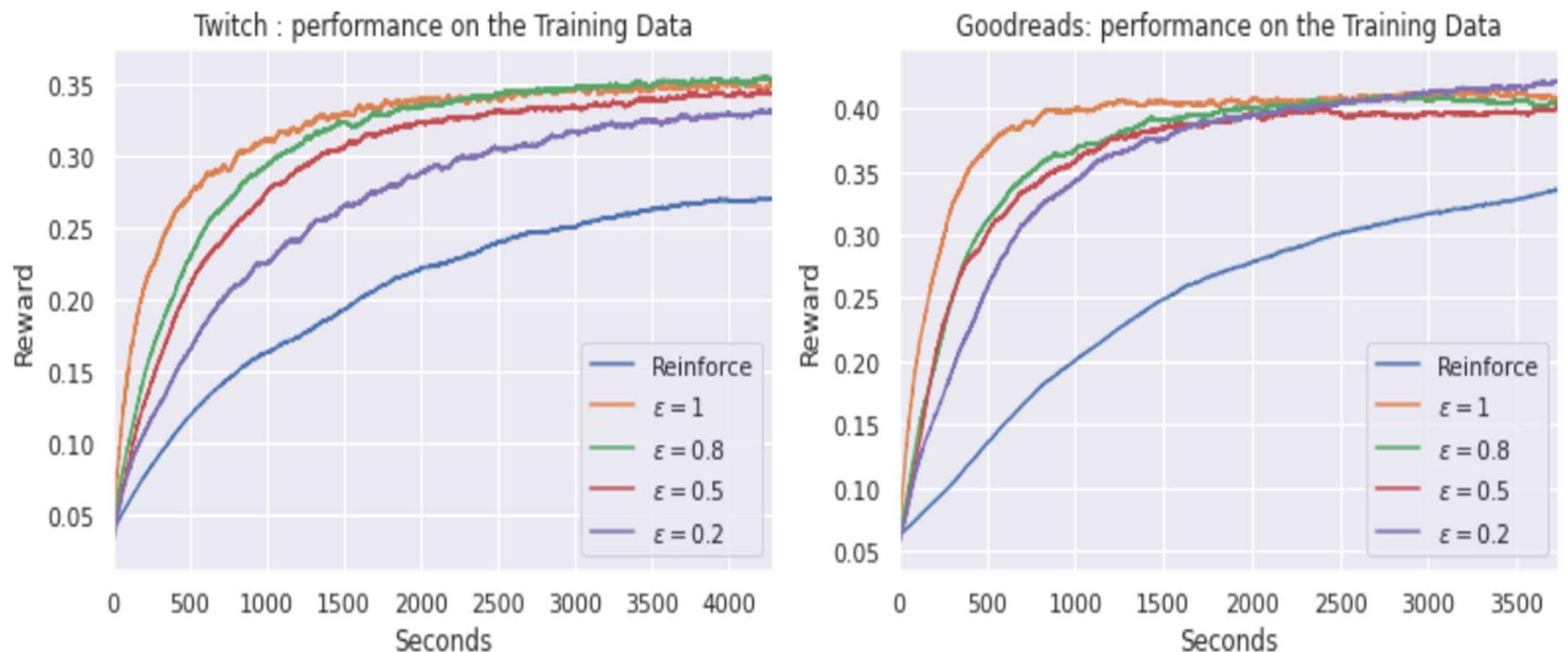


Fig. 5. Performance plots given time : ϵ 's effect on the training

[Sakhi et al. 2022](#)

Notebook 2

[https://colab.research.google.com/github/otmhi/Reward-Optimizing-Reco/blob/main/Reward Optimizing Slate Recommendation with DL and MIPS Part 2.ipynb](https://colab.research.google.com/github/otmhi/Reward-Optimizing-Reco/blob/main/Reward%20Optimizing%20Slate%20Recommendation%20with%20DL%20and%20MIPS%20Part%202.ipynb)

Tutorial covering Horvitz Thompson for Recommender Systems

[RecSys2021 Tutorial](#)